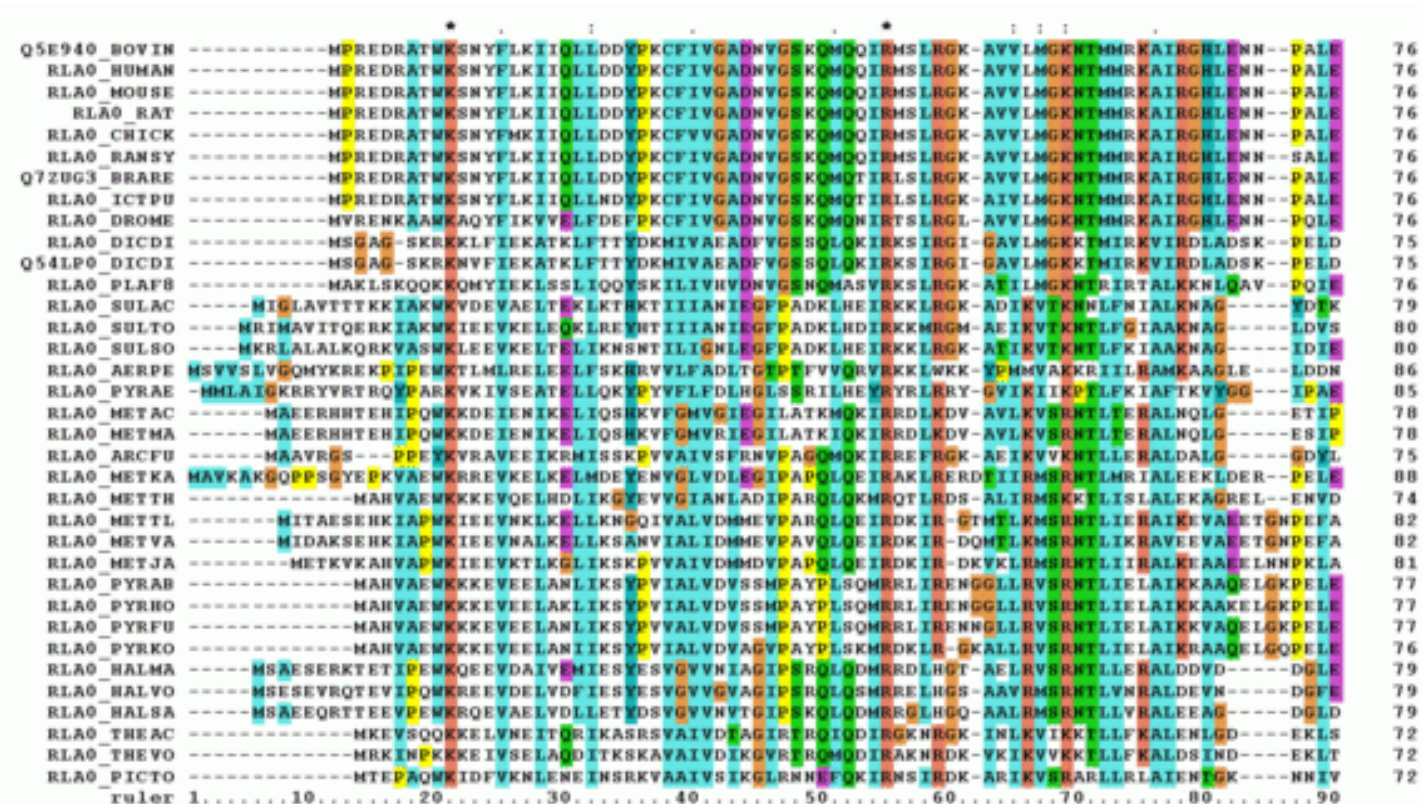


EECS730: Introduction to Bioinformatics

Lecture 06: Multiple Sequence Alignment



https://upload.wikimedia.org/wikipedia/commons/thumb/7/79/RPLP0_90_ClustalW_aln.gif/575px-RPLP0_90_ClustalW_aln.gif

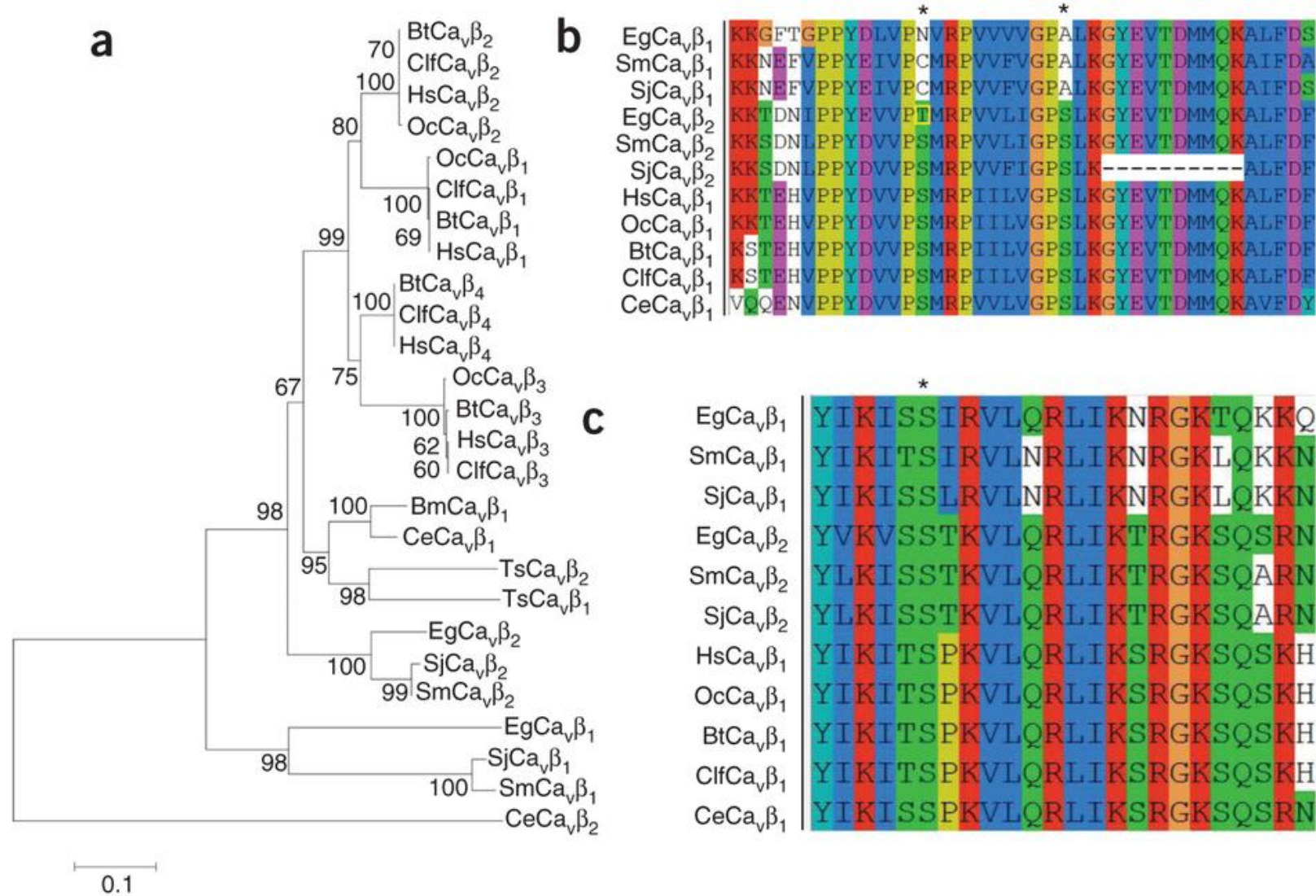
Slides adapted from Dr. Shaojie Zhang (University of Central Florida)

Multiple alignments

- Reveal evolutionary history (speciation-related mutations)
- Prediction of protein structure and protein function
- Determine consensus sequence for sequence assembly

- Generalization of the pairwise alignment algorithm

Multiple alignments



Object function

- To maximize the conservation of the alignment columns
- The more conserved the columns, the better the alignment
- Three scoring functions to characterize the conservation of the columns
 - Multiple Longest Common Sequence
 - Entropy
 - Sum-of-pair scores

Multiple Longest Common Subsequence

- A column is a “match” if all the letters in the column are the same

A
A
A
A
A
A
A
A
A
A

- Similar idea to the LCS problem formulation for pairwise alignment
- Only good for very similar sequences

Entropy

- Define frequencies for the occurrence of each letter in each column of multiple alignment (gap may be included into the alphabet)
 - $p_A = 1, p_T=p_G=p_C=0$ (1st column)
 - $p_A = 0.75, p_T = 0.25, p_G=p_C=0$ (2nd column)
 - $p_A = 0.50, p_T = 0.25, p_C=0.25, p_G=0$ (3rd column)
- Compute entropy of each column

$$- \sum_{X=A,T,G,C} p_X \log p_X$$

AAA
AAA
AAT
ATC

Entropy cont.

- Best case

$$\text{entropy} \begin{pmatrix} A \\ A \\ A \\ A \end{pmatrix} = 0$$

- Worst case

$$\text{entropy} \begin{pmatrix} A \\ T \\ G \\ C \end{pmatrix} = -\sum \frac{1}{4} \log \frac{1}{4} = -4 \left(\frac{1}{4} * -2 \right) = 2$$

- Entropy for a multiple alignment is the sum of entropies of its columns

Sum-of-pair score

Every multiple alignment induces pairwise alignments

x: AC-GCGG-C
y: AC-GC-GAG
z: GCCGC-GAG

Induces:

x: ACGCGG-C; **x**: AC-GCGG-C; **y**: AC-GCGAG
y: ACGC-GAC; **z**: GCCGC-GAG; **z**: GCCGCAG

Sum-of-pair score cont.

- The alignment score for the multiple alignment is the sum of the alignment scores of all of its induced pairwise alignments

- Consider pairwise alignment of sequences

a_i and a_j

imposed by a multiple alignment of k sequences

- Denote the score of this suboptimal (not necessarily optimal) pairwise alignment as

$$s^*(a_i, a_j)$$

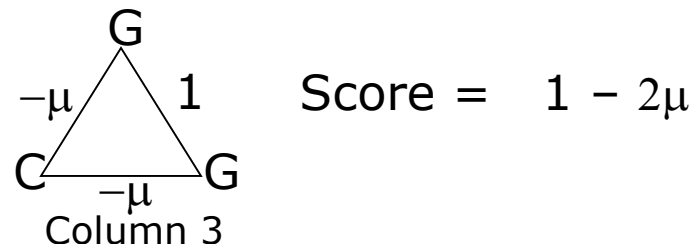
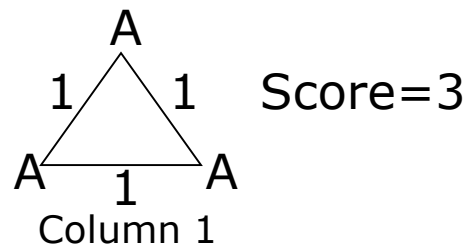
- Sum up the pairwise scores for a multiple alignment:

$$s(a_1, \dots, a_k) = \sum_{i,j} s^*(a_i, a_j)$$

Sum-of-pair score cont.

- It can also be computed column-wise
- This is useful for dynamic programming algorithm that breaks the problem into smaller sub-problems

a_1 ATG-C-AAT
· A-G-CATAT
 a_k ATCCCATTT



How to compute optimal multiple alignment

- Extending the dynamic programming algorithm for pairwise alignment
- Recall what does the score mean for each entry in the 2D dynamic programming table for pairwise alignments
- What is the dimension of the multiple sequence alignment dynamic programming table and what should we store there?

How to compute optimal multiple alignment

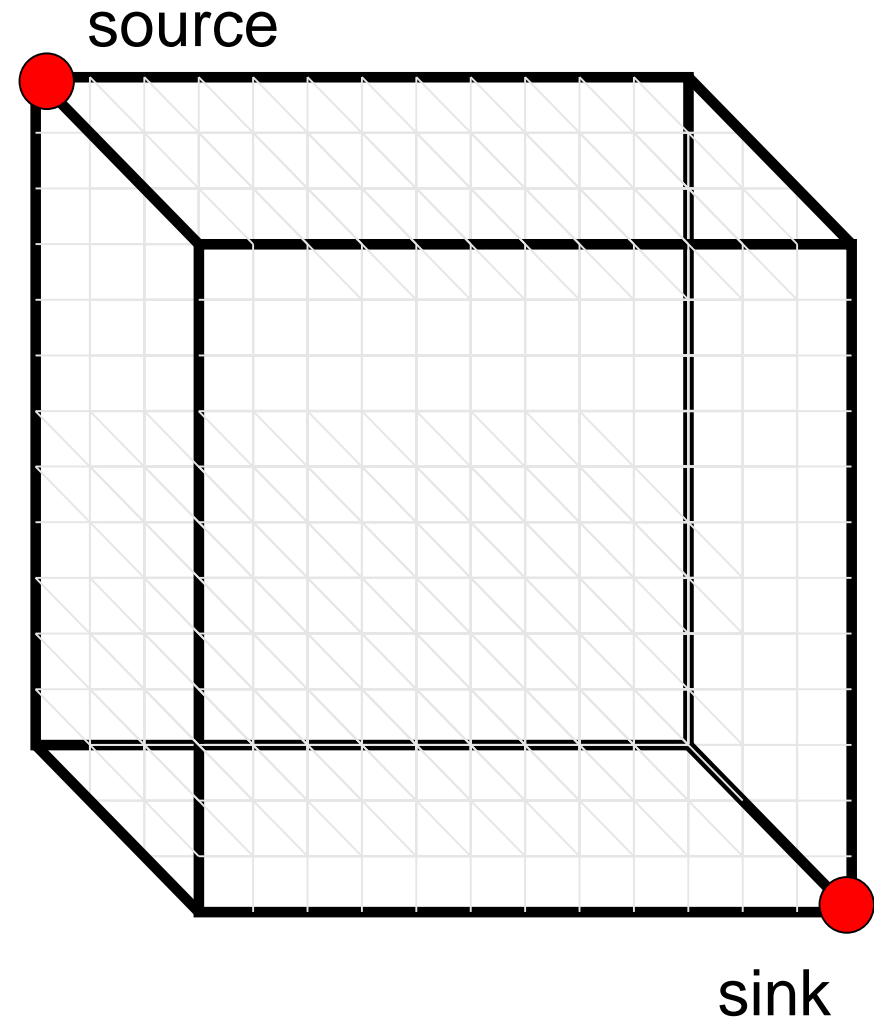
- Each entry in the 2D DP table stores the best score for aligning the prefixes of the two sequences:
 - The entry (i, j) stores alignment score between $S1(0, i)$ and $S2(0, j)$, where $S1$ and $S2$ are the two sequences being aligned.
- This can also be extended to multiple alignment case
- How many different combinations of prefixes alignment for n sequences?
 - $l_1 * l_2 * \dots * l_n$, where l is the length of a given sequence
 - So the DP table for multiple alignment is an n -dimensional table
 - It degenerates to 2D table for pairwise alignment

How to compute optimal multiple alignment

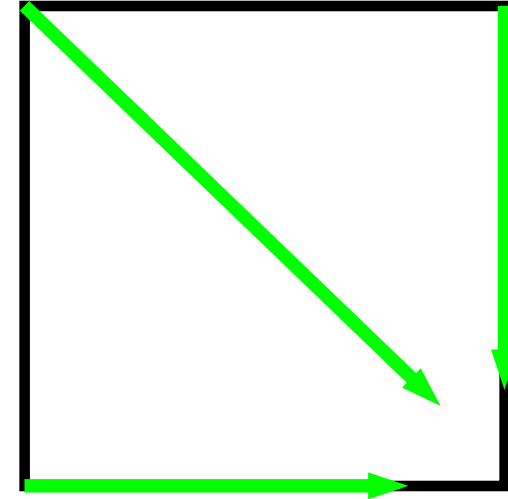
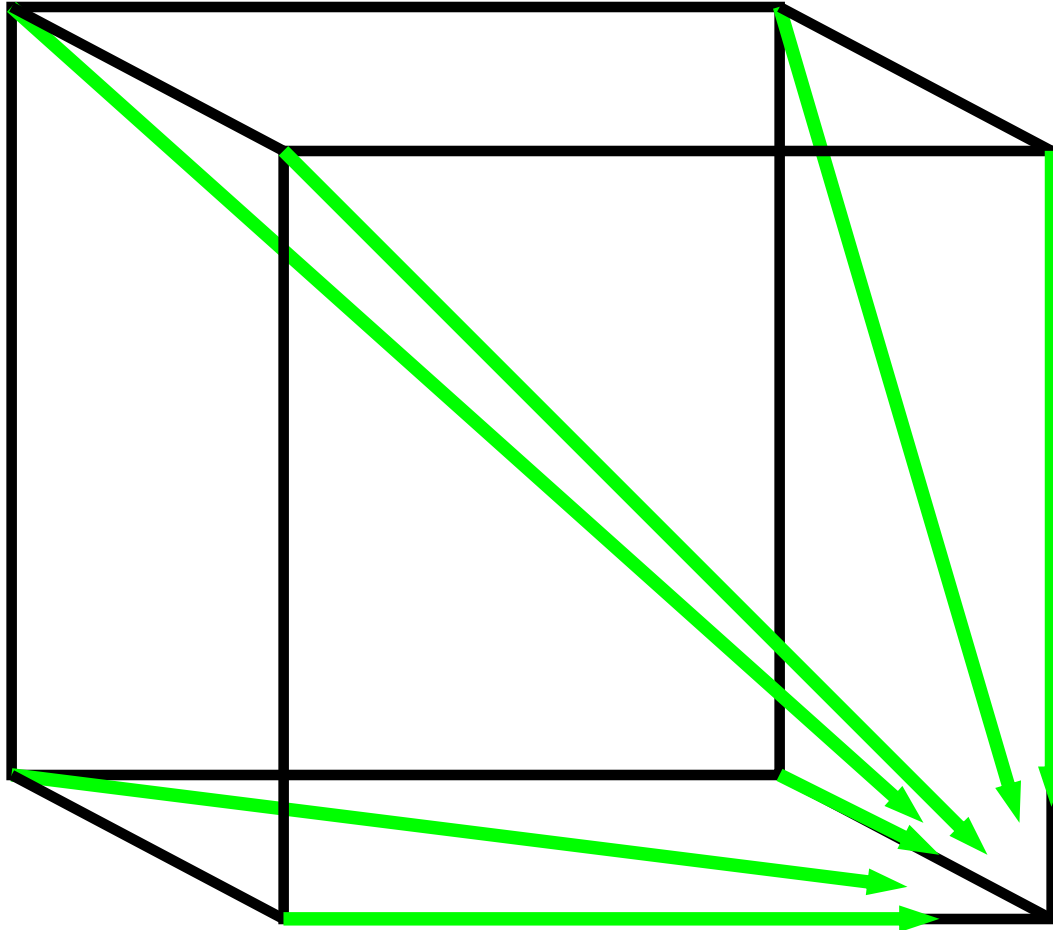
- Now, how many entries do we need to refer to in order to compute the score of an entry?
- Recall that each stage of the DP algorithm append one or zero character of each sequence to the existing alignment
- There are two choices (one or zero character), and there are n sequences in total, so there are 2^n entries to refer in total
- More precisely, we do not allow a column of all gaps, which means the combination of all zeros is invalid, and it reduces the number of entries to refer to as $2^n - 1$
- For pairwise alignments, we need to refer to $2^2 - 1 = 3$ entries in the DP table (left, up, and upper-left)

An example for aligning three sequences

- A three-dimensional Manhattan Tourist Problem
- The DP matrix is 3D
- We aim at finding the path that corresponds to the best alignment



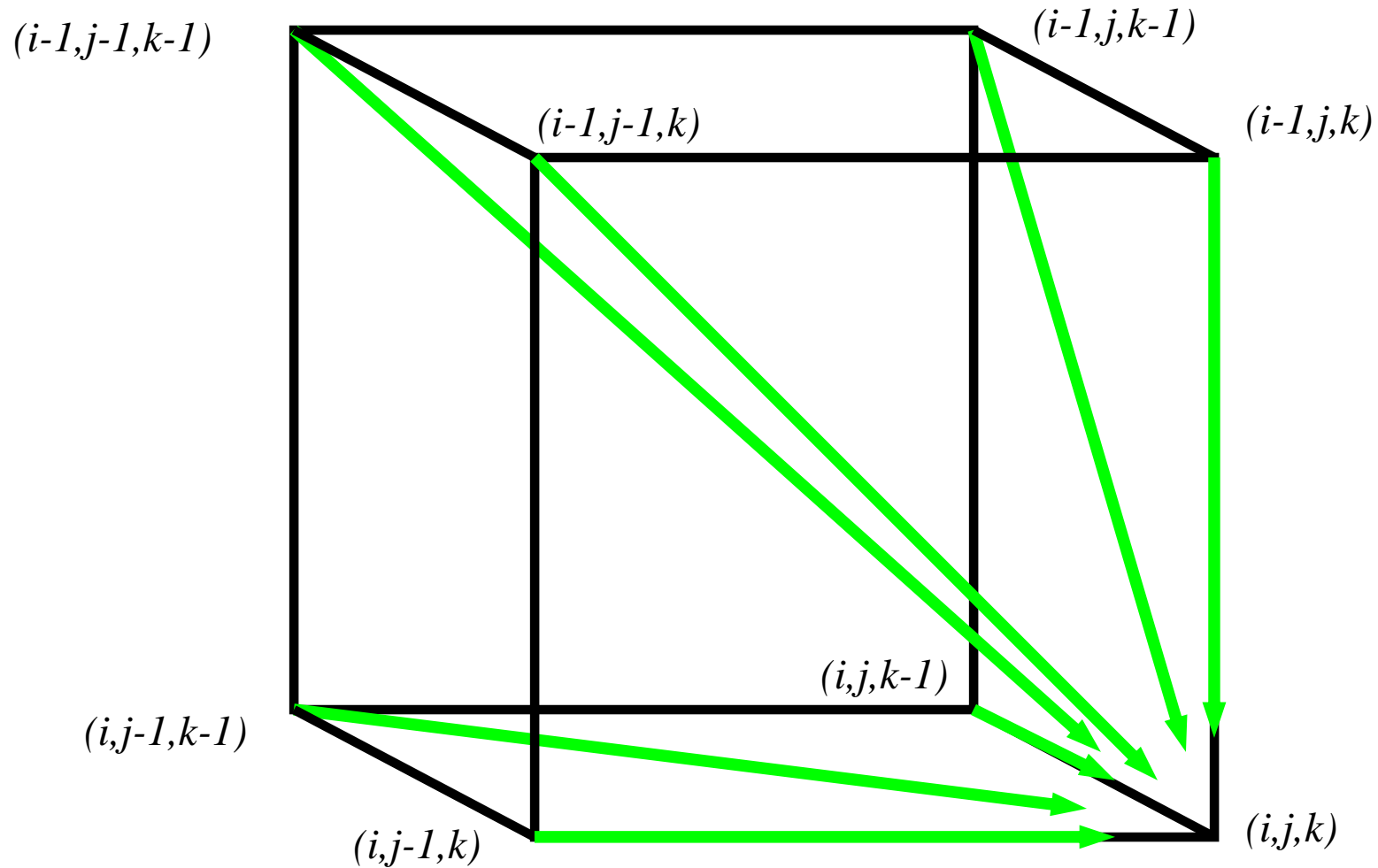
Filling the DP matrix



In **2-D**, 3 edges
in each unit
square

In **3-D**, 7 edges
in each unit cube

Architecture of the 3D alignment cell



Recursive function for MSA

$$\bullet s_{i,j,k} = \max \left\{ \begin{array}{l}
 s_{i-1,j-1,k-1} + \delta(v_i, w_j, u_k) \\
 s_{i-1,j-1,k} + \delta(v_i, w_j, _) \\
 s_{i-1,j,k-1} + \delta(v_i, _, u_k) \\
 s_{i,j-1,k-1} + \delta(_, w_j, u_k) \\
 s_{i-1,j,k} + \delta(v_i, _, _) \\
 s_{i,j-1,k} + \delta(_, w_j, _) \\
 s_{i,j,k-1} + \delta(_, _, u_k)
 \end{array} \right.$$

cube diagonal:
no indels

 face diagonal:
one indel

 edge diagonal:
two indels

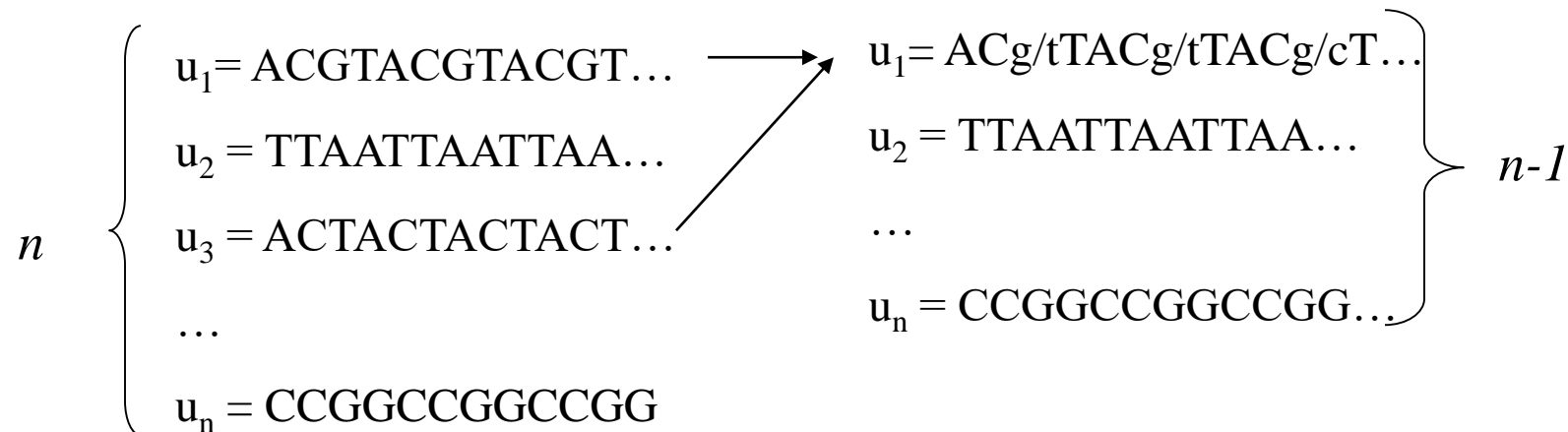
- $S(x, y, z)$ is an entry in the 3-D scoring matrix
- $\delta(x, y, z)$ can be computed as the sum-of-pair score

What is the time complexity?

- We have l^n entries to fill, filling each entry takes 2^n time
- The overall complexity is $O(l^n * 2^n)$
- Conclusion: dynamic programming approach for alignment between two sequences is easily extended to n sequences but it is impractical due to exponential running time.

Progressive multiple sequence alignment

- Perform all-against-all pairwise alignments for the n sequences
- Choose most similar pair of strings and combine into a profile , thereby reducing alignment of n sequences to an alignment of $n-1$ sequences/profiles. **Repeat until 1 sequence/profile remains**
- This is a heuristic greedy method



Completing the iteration

- We need to find the most similar pair of strings at each iteration
- Therefore we need to redefine the similarity between the newly summarized profile and the other strings/profiles
- Using the Neighbor Joining Algorithm

$$d(u, k) = \frac{1}{2} [d(f, k) + d(g, k) - d(f, g)]$$

https://en.wikipedia.org/wiki/Neighbor_joining

- Here, u is the new profile, and f and g are the two strings/profiles that form u , and k is an arbitrary string/profile remains

Completing the iteration

- Profile representation

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

- A single sequence can be viewed as a special case profile

Aligning two profiles

- Two profiles represented using frequencies can be aligned using slightly modified pairwise sequence alignment algorithm
- The score for matching two columns can be computed as the sum-of-pair scores of the two columns
- Affine gap penalty can be easily incorporated

Example

DISTANCES between protein sequences:

Calculated over: 1 to 167

Correction method: Simple distance (no corrections)

Distances are: observed number of substitutions per 100 amino acids

Symmatrix version 1

Number of matrices: 1

//

Matrix 1, dimension: 7

Key for column and row indices:

1 hba_human
2 hba_horse
3 hbb_human
4 hbb_horse
5 glb5_petma
6 myg_phyca
7 lgb2_luplu

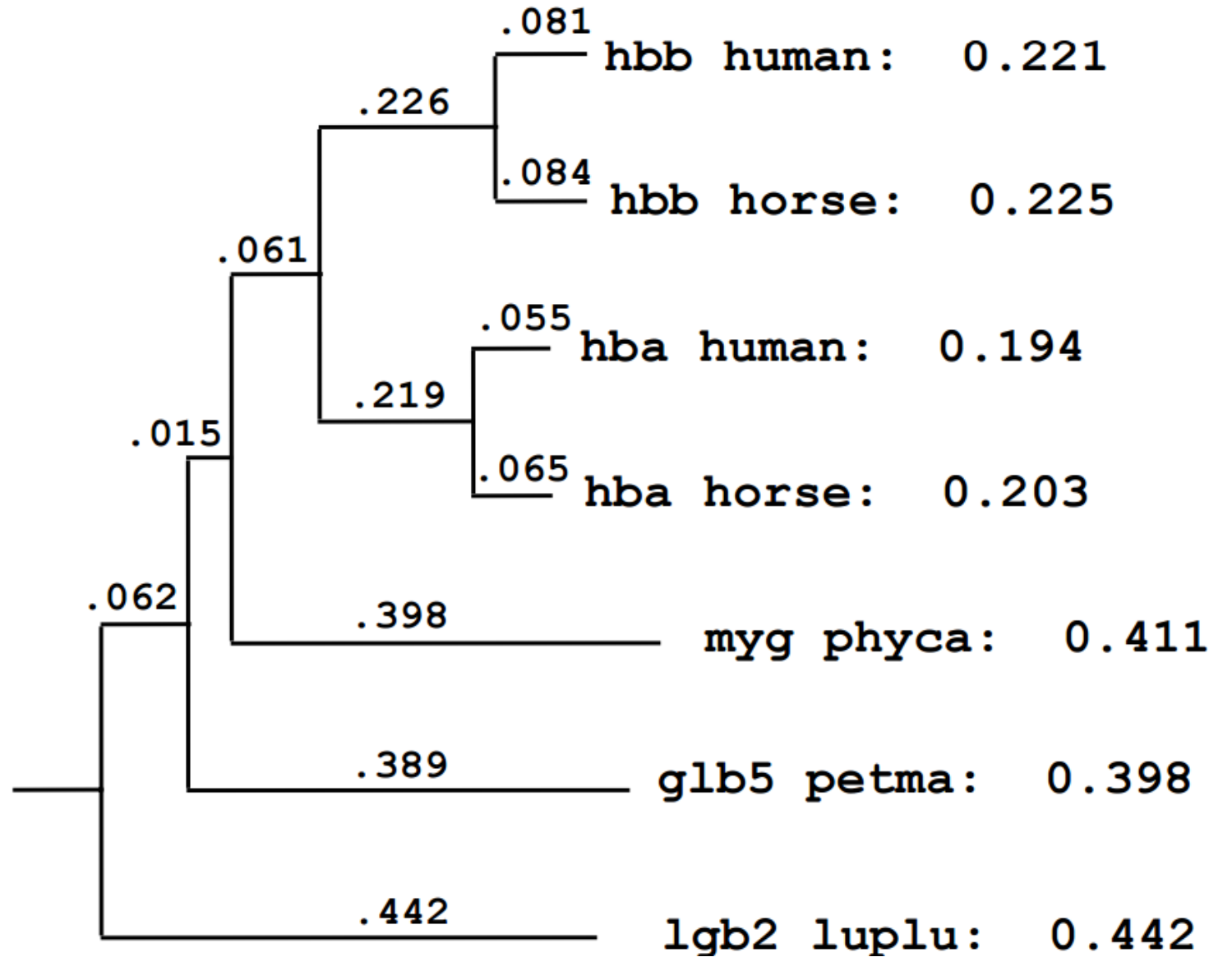
Matrix 1: Part 1

	1	2	3	4	5	6	7	..
1	0.00	12.06	54.68	55.40	64.12	71.74	83.57	
2		0.00	55.40	53.96	64.89	72.46	82.86	
3			0.00	16.44	74.26	73.94	82.52	
4				0.00	75.74	73.94	81.12	
5					0.00	75.91	82.61	
6						0.00	80.95	
7							0.00	

The guide tree

For computation of the distances see

https://en.wikipedia.org/wiki/Neighbor_joining



Progressive alignment based on the guide tree

In our globin example, we align in the following order:

- a) human and horse beta-globin;
- b) human and horse alpha-globin;
- c) the two beta-globins and the two alpha-globins;
- d) myoglobin and the haemoglobins;
- e) cyanohaemoglobin and the combined haemoglobin, myoglobin group;
- f) leghaemoglobin and the rest.

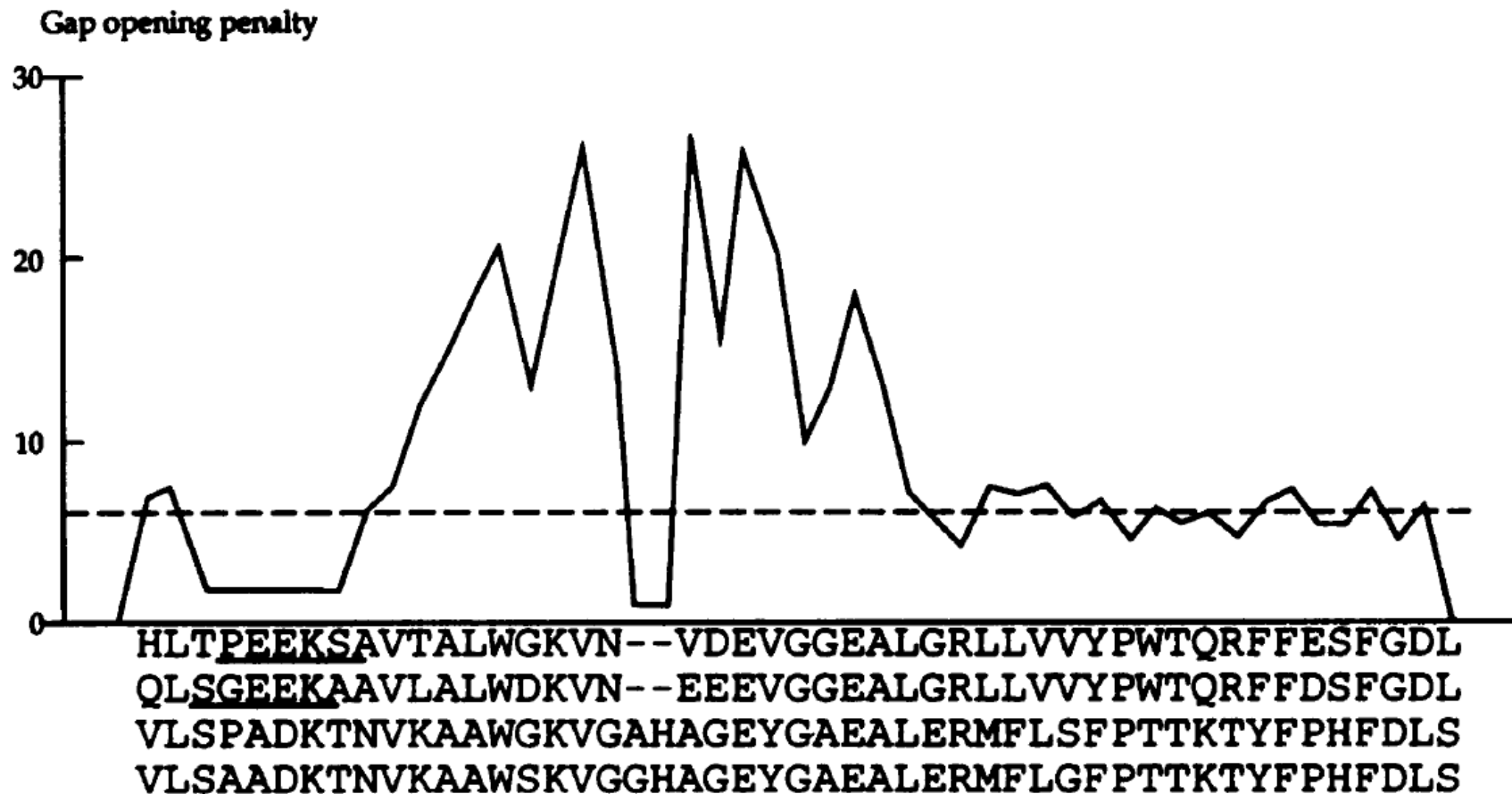
Time complexity for progressive multiple alignment

- N^3
- At each stage, we need to re-compute the distance between the newly formed profile and the other N (in the worst case) sequences/profiles (linear)
- At each stage, we also need to perform pairwise alignment (square)
- Taken together, each stage requires square time
- We have N stages because each stage we reduce the size of set of sequences/profiles by 1
- So N^3 in total

ClustalW: more sophisticated scoring function

- Firstly, individual weights are assigned to each sequence in a partial alignment in order to downweight near-duplicate sequences and up-weight the most divergent ones.
- Secondly, amino acid substitution matrices are varied at different alignment stages according to the divergence of the sequences to be aligned.
- Thirdly, residue-specific gap penalties and locally reduced gap penalties in hydrophilic regions encourage new gaps in potential loop regions rather than regular secondary structure.
- Fourthly, positions in early alignments where gaps have been opened receive locally reduced gap penalties to encourage the opening up of new gaps at these positions.

ClustalW: residue-dependent gap penalty



ClustalW: output

```

                10         20         30         40         50         60
Hbb_Human.pep  -----VHLTPEEKSAVTALWGKVN--VDEVGGEALGRLLVYYPWTCRFFESFGDLST
Hbb_Horse.pep  -----VQLSGEEKAAVLALWDKVN--EEEVGGGEALGRLLVYYPWTCRFFDSFGDLSN
Hba_Human.pep  -----VLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPFDLS--
Hba_Horse.pep  -----VLSAADKTNVKAAWSKVGGHAGEYGAEALERMF LGFPPTTKTYFPFDLS--
Myg_Phyca.pep  -----VLSSEGEWQLVLHVWAKVEADVAGHGQDILIRLFKSHPETLEKFD RFKHLKT
Glb5_Petma.pep PIVDTGSVAPLSAAEKT KIRSAWAPVYSTYETSGVDILVKFFTSTPAAC EFPKFKGLTT
Lgb2_Luplu.pep -----GALTESQAALVKSSWEEFNANIPKHTHRFFILVLEIAPA AKDLFSFLKGTSE
                *         *         *         *

```

```

Hbb_Human.pep  PDAVMGNPKVKAHGKKV LGA FSDGLAHL D-----NLKGT FATLSELHC DKLHVDPENFRL
Hbb_Horse.pep  PGAVMGNPKVKAHGKKV LHSFGE GVHHL D-----NLKGT FAALSELHC DKLHVDPENFRL
Hba_Human.pep  ----HGSAQVKGHGKKV ADAL TNAVAHVD-----DMPNALSALSDLHA HKLRVDPVNFKL
Hba_Horse.pep  ----HGSAQVKAHGKKV GDAL TLAVGHLD-----DLPGALS NLSDLHA HKLRVDPVNFKL
Myg_Phyca.pep  EAEMKASEDLKKHGVT VLTALGAILKKKG-----HHEAELKPLAQSHATKHKIPIKYLEF
Glb5_Petma.pep ADQLKKSADVRWHAERI INAVND AVASMDDT--EKMSMKLRDL SGKHAKSFQVDPQYFKV
Lgb2_Luplu.pep VP--QNNPELQAHAGKV FKL VYEAAIQLOVTGVVVT DATLKNLGSVHVSKG-VADAHFPV
                .. *         *

```

```

Hbb_Human.pep  LGNVLVCVLAH HFGKEFT PPVQAAYQKV VAGVANALAHKYH-----
Hbb_Horse.pep  LGNVLVVVLAH HFGKDFTP ELQASYQKV VAGVANALAHKYH-----
Hba_Human.pep  LSHCLLVTLAA HLP AEF TPAVHASLDKFLASVSTVLT SKYR-----
Hba_Horse.pep  LSHCLLSTLAV HLPNDFT PAVHASLDKFLSSVSTVLT SKYR-----
Myg_Phyca.pep  ISEAI IHVLSRHPGDFGADAQGAMNKALELFRKDIAAKYKELGYQG
Glb5_Petma.pep LAAVIADTVAAG-----DAGFEKLMSMICILLRSAY-----
Lgb2_Luplu.pep VKEAILKTIKEVVGAKWSEELNSAWTIAYDELAIVIKKEMNDAA---

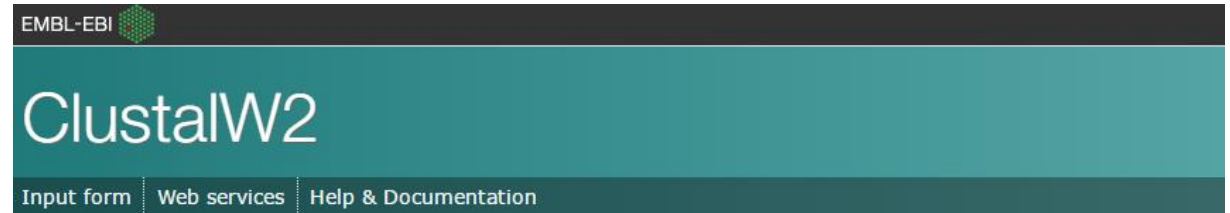
```

7 α -helices

ClustalW: multiple sequence alignment

- <http://www.ebi.ac.uk/Tools/msa/clustalw2/>

- <http://www.genome.jp/tools/clustalw/>



Multiple Sequence Alignment by CLUSTALW

ETE3	MAFFT	CLUSTALW	PRRN
Help			
General Setting Parameters:			
Output Format: <input type="text" value="CLUSTAL"/>			
Pairwise Alignment: <input checked="" type="radio"/> FAST/APPROXIMATE <input type="radio"/> SLOW/ACCURATE			
Enter your sequences (with labels) below (copy & paste): <input checked="" type="radio"/> PROTEIN <input type="radio"/> DNA			
Support Formats: FASTA (Pearson), NBRF/PIR, EMBL/Swiss Prot, GDE, CLUSTAL, and GCG/MSF			
<input type="text"/>			
Or give the file name containing your query			
<input type="button" value="Choose File"/> No file chosen			
<input type="button" value="Execute Multiple Alignment"/> <input type="button" value="Reset"/>			