

## Multiple Sequence Alignment

		*	:	:	*	:	:	:	ruler
Q5E940_BOVIN	-----	M P R E D R A T W K S N Y F L K I I Q L L D D Y P K C F I V G A D N V G S K Q M Q Q I R M S L R G K - A V V L M G K N T M M R K A I R G H L E N N -- P A L E							76
RLA0_HUMAN	-----	M P R E D R A T W K S N Y F L K I I Q L L D D Y P K C F I V G A D N V G S K Q M Q Q I R M S L R G K - A V V L M G K N T M M R K A I R G H L E N N -- P A L E							76
RLA0_MOUSE	-----	M P R E D R A T W K S N Y F L K I I Q L L D D Y P K C F I V G A D N V G S K Q M Q Q I R M S L R G K - A V V L M G K N T M M R K A I R G H L E N N -- P A L E							76
RLA0_RAT	-----	M P R E D R A T W K S N Y F L K I I Q L L D D Y P K C F I V G A D N V G S K Q M Q Q I R M S L R G K - A V V L M G K N T M M R K A I R G H L E N N -- P A L E							76
RLA0_CHICK	-----	M P R E D R A T W K S N Y F M K I I Q L L D D Y P K C F V V G A D N V G S K Q M Q Q I R M S L R G K - A V V L M G K N T M M R K A I R G H L E N N -- P A L E							76
RLA0_RANSY	-----	M P R E D R A T W K S N Y F L K I I Q L L D D Y P K C F I V G A D N V G S K Q M Q Q I R M S L R G K - A V V L M G K N T M M R K A I R G H L E N N -- S A L E							76
Q7ZUG3_BRARE	-----	M P R E D R A T W K S N Y F L K I I Q L L D D Y P K C F I V G A D N V G S K Q M Q T I R L S L R G K - A V V L M G K N T M M R K A I R G H L E N N -- P A L E							76
RLA0_ICTPU	-----	M P R E D R A T W K S N Y F L K I I Q L L N D Y P K C F I V G A D N V G S K Q M Q T I R L S L R G K - A I V L M G K N T M M R K A I R G H L E N N -- P A L E							76
RLA0_DROME	-----	M V R E N K A A W K A Q Y F I K V V E L F D E F P K C F I V G A D N V G S K O M O N I R T S L R G L - A V V L M G K N T M M R K A I R G H L E N N -- P Q L E							76
RLA0_DICDI	-----	M S G A G - S K R K K L F I E K A T K L F T T Y D K M I V A E A D F V G S S Q L Q K I R K S I R G I - G A V L M G K K T M I R K V I R D L A D S K -- P E L D							75
Q54LP0_DICDI	-----	M S G A G - S K R K N V F I E K A T K L F T T Y D K M I V A E A D F V G S S Q L Q K I R K S I R G I - G A V L M G K K T M I R K V I R D L A D S K -- P E L D							75
RLA0_PLAF8	-----	M A K L S K Q Q K K Q M Y I E K L S S I Q Q Y S K I L I V H V D N V G S N Q M A S V R K S L R G K - A T I L M G K N T R I R T A L K K N L Q A V -- P Q I E							76
RLA0_SULAC	-----	M I G L A V T T T K K I A K W K V D E V A E L T E K L K T H K T I I I A N I E G F P A D K L H E I R K K L R G K - A D I K V T K N N L F N I A L K N A G -- Y D T K							79
RLA0_SULTO	-----	M R I M A V I T Q E R K I A K W K I E E V K E L E Q K L R E Y H T I I I A N I E G F P A D K L H D I R K K M R G M - A E I K V T K N T L F G I A A K N A G -- L D V S							80
RLA0_SULSO	-----	M K R L A L A L K Q R K V A S W K L E E V K E L T E L I K N S N T I L I G N L E G F P A D K L H E I R K K L R G K - A T I K V T K N T L F K I A A K N A G -- I D I E							80
RLA0_AERPE	M S V V S L V G Q M Y K R E K P I P E W K T I M L R E L E F S K H R V V L F A D L T G T P T F V V Q R V R K K L W K K - Y P M M V A K K R I I L R A M K A A G L E -- L D D N								86
RLA0_PYRAE	-M M L A I G K R R Y V R T R Q Y P A R K V K I V S E A T E L L Q K Y P V V F L D H G L S S R I L H E Y R Y R L R R Y - G V I K I I K P T L F K I A F T K V Y G G -- I P A E								85
RLA0_METAC	-----	M A E E R H H T E H I P Q W K K D E I E N I K E L I Q S H K V F G M V G I E G I L A T K M Q K I R R D L K D V - A V L K V S R N T L T E R A L N Q L G -- E T I P							78
RLA0_METMA	-----	M A E E R H H T E H I P Q W K K D E I E N I K E L I Q S H K V F G M V R I E G I L A T K I Q K I R R D L K D V - A V L K V S R N T L T E R A L N Q L G -- E S I P							78
RLA0_ARCFU	-----	M A A V R G S -- P P E Y K V R A V E E I K R M I S S K P V V A I V S F R N V P A G Q M Q K I R R E F R G K - A E I K V V K N T L L E R A L D A L G -- G D Y L							75
RLA0_METKA	M A V K A K G Q P P S G Y E P K V A E W K R R E V K E L K E L M D E Y E N V G L V D L E G I P A P Q L Q E I R A K L R E R D T I I R M S R N T L M R I A L E E K L D E R -- P E L E								88
RLA0_METTH	-----	M A H V A E W K K K E V Q E L H D L I K G Y E V V G I A N L A D I P A R Q L Q K M R Q T L R D S - A L I R M S K K T L I S L A L E K A G R E L -- E N V D							74
RLA0_METTL	-----	M I T A E S E H K I A P W K I E E V N K L K E L L K N G Q I V A L V D M M E V P A R Q L Q E I R D K I R - G T M T L K M S R N T L I E R A I K E V A E E T G N P E F A							82
RLA0_METVA	-----	M I D A K S E H K I A P W K I E E V N A L K E L L K S A N V I A L I D M M E V P A V Q L Q E I R D K I R - D Q M T L K M S R N T L I K R A V E E V A E E T G N P E F A							82
RLA0_METJA	-----	M E T K V A H V A P W K I E E V K T L K G L I K S K P V V A I V D M M D V P A P Q L Q E I R D K I R - D K V K L R M S R N T L I I R A L K E A A E E L N N P K L A							81
RLA0_PYRAB	-----	M A H V A E W K K K E V E E L A N L I K S Y P V I A L V D V S S M P A Y P L S Q M R R L I R E N G G L L R V S R N T L I E L A I K K A A Q E L G K P E L E							77
RLA0_PYRHO	-----	M A H V A E W K K K E V E E L A K L I K S Y P V I A L V D V S S M P A Y P L S Q M R R L I R E N G G L L R V S R N T L I E L A I K K A A K E L G K P E L E							77
RLA0_PYRFU	-----	M A H V A E W K K K E V E E L A N L I K S Y P V V A L V D V S S M P A Y P L S Q M R R L I R E N G G L L R V S R N T L I E L A I K K V A Q E L G K P E L E							77
RLA0_PYRKO	-----	M A H V A E W K K K E V E E L A N I I K S Y P V I A L V D V A G V P A Y P L S K M R D K L R - G K A L L R V S R N T L I E L A I K R A A Q E L G Q P E L E							76
RLA0_HALMA	M S A E S E R K T E T I P E W K Q E E V D A I V E M I E S Y E S V G V V N I A G I P S R Q L Q D M R R D L H G T - A E L R V S R N T L L E R A L D D V D -- D G L E								79
RLA0_HALVO	M S E S E V R Q T E V I P Q W K R E E V D E L V D F I E S Y E S V G V V V G V A G I P S R Q L Q S M R R E L H G S - A A V R M S R N T L V N R A L D E V N -- D G F E								79
RLA0_HALSA	M S A E E Q R T T E E V P E W K R Q E V A E L V D L L E T Y D S V G V V V N V T G I P S K Q L Q D M R R G L H G Q - A A L R M S R N T L L V R A L E E A G -- D G L D								79
RLA0_THEAC	-----	M K E V S Q Q K K E L V N E I T Q R I K A S R S V A I V D T A G I R T R Q I Q D I R G K N R G K - I N L K V I K K T L L F K A L E N L G D -- E K L S							72
RLA0_THEVO	-----	M R K I N P K K K E I V S E L A Q D I T K S K A V A I V D I K G V R T R Q M Q D I R A K N R D K - V K I K V V K K T L L F K A L D S I N D -- E K L T							72
RLA0_PICTO	-----	M T E P A Q W K I D F V K N L E N E I N S R K V A A I V S I K G L R N N F E Q K I R N S I R D K - A R I K V S R A R L L R A I E N T G K -- N N I V							72

ruler 1.....10.....20.....30.....40.....50.....60.....70.....80.....90

# Overview

## Multiple Sequence Alignment

- introduction
- methods for Multiple Sequence Alignment (MSA)
- uses of MSA

## Learning Objectives

- describe the uses of MSA
- explain how ClustalW performs MSA
- describe alternative methods for MSA (Muscle, ProbCons)
- understand how to interpret MSA results

# Multiple Sequence Alignment (MSA)

- A multiple sequence alignment is a collection of  $\geq 3$  nucleotide or protein sequences that are aligned
- These alignments can reveal domains of highly related sequence that in turn can be used to define regions of interest
- For proteins that could be “domains” that define a protein family
- For nucleotides that could be conserved regions that define coding, regulatory or structural elements of a genome
- The assumption is that sequences that share similarity share function i.e. they are homologous
- MSAs can be used to “test” whether a novel sequence is a member of a family, contains a particular domain, or has some regulatory function in common with known regulators for example

# Homology (Reminder)

**Orthologs** are the result of a speciation event  
**Paralogs** are the result of a duplication event

All are Homologs

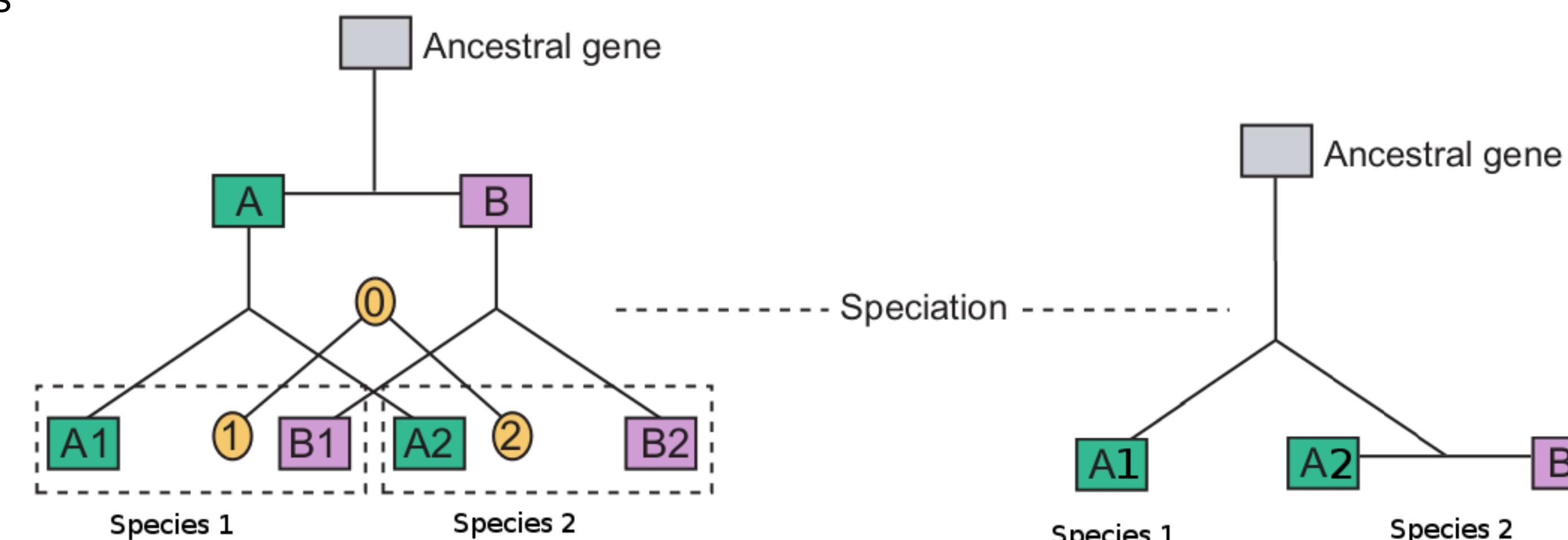
Paralogs

A1 and B1  
A1 and B2

A2 and B1  
A2 and B2

Orthologs

A1 and A2  
B1 and B2



All are Homologs

Paralogs

A2 and B2 are paralogs of each other

Orthologs

A2 and B2 are orthologs of A1

# Multiple Sequence Alignment (MSA)

MSAs are easy to create for closely related sequences such as glyceraldehyde-3-phosphate dehydrogenase (GAPDH)

<a href="#">NP_002037.2</a>	194	KLWRDGRGALQNIIPASTGAAKAVGVIPELNGKLTGMAFRVPTANVSVV	243
<a href="#">XP_508955.1</a>	194	KLWRDGRGALQNIIPASTGAAKAVGVIPELNGKLTGMAFRVPTANVSVV	243
<a href="#">XP_001105471.1</a>	194	KLWRDGRGALQNIIPASTGAAKAVGVIPELNGKLTGMAFRVPTANVSVV	243
<a href="#">NP_001003142.1</a>	192	KMWRDGRGAAQNIIPASTGAAKAVGVIPELNGKLTGMAFRVPTPNVSVV	241
<a href="#">XP_003435697.1</a>	192	KLWRDGRGAAQNIIPASTGAAKAVGVIPELNGKLTGMAFRVPTPNVSVV	241
<a href="#">XP_003434435.1</a>	192	KLWRDGRGAAQNIIPASTGAAKAVGVIPELNGKLTGMAFCVPTPNVSVV	241
<a href="#">NP_001029206.1</a>	192	KLWRDGRGAAQNIIPASTGAAKAVGVIPELNGKLTGMAFRVPTPNVSVV	241
<a href="#">NP_032110.1</a>	192	KLWRDGRGAAQNIIPASTGAAKAVGVIPELNGKLTGMAFRVPTPNVSVV	241
<a href="#">XP_001476757.1</a>	192	KLWRDGRGAAQNIIPASTGAAKAVGVIPELNGKLTGMAFRVPTPNVSVV	241
<a href="#">NP_058704.1</a>	192	KLWRDGRGAAQNIIPASTGAAKAVGVIPELNGKLTGMAFRVPTPNVSVV	241
<a href="#">NP_989636.1</a>	192	KLWRDGRGAAQNIIPASTGAAKAVGVIPELNGKLTGMAFRVPTPNVSVV	241
<a href="#">NP_001108586.1</a>	192	KLWRDGRGASQNIIPASTGAAKAVGVIPELNGKLTGMAFRVPTPNVSVV	241
<a href="#">NP_001259584.1</a>	191	KLWRDGRGAAQNIIPASTGAAKAVGVIPALNGKLTGMAFRVPTPNVSVV	240
<a href="#">NP_525108.2</a>	191	KLWRDGRGAAQNIIPAATGAAKAVGVIPALNGKLTGMAFRVPTPNVSVV	240
<a href="#">XP_318655.2</a>	191	KLWRDGRGAAQNIIPAATGAAKAVGVIPALNGKLTGMAFRVPTPNVSVV	240
<a href="#">NP_496237.1</a>	200	KLWRDGRGAGQNIIPASTGAAKAVGVIPELNGKLTGMAFRVPTPDVSVV	249
<a href="#">NP_496192.1</a>	200	KLWRDGRGAGQNIIPASTGAAKAVGVIPELNGKLTGMAFRVPTPDVSVV	249
<a href="#">NP_508534.3</a>	200	KLWRDGRGAGQNIIPASTGAAKAVGVIPELNGKLTGMAFRVPTPDVSVV	249
<a href="#">NP_508535.1</a>	200	KLWRDGRGAGQNIIPASTGAAKAVGVIPELNGKLTGMAFRVPTPDVSVV	249
<a href="#">NP_012483.3</a>	192	KDWRGGRTASNIIIPSSTGAAKAVGVKVLPELQGKLTGMAFRVPTDVSVV	241
<a href="#">NP_011708.3</a>	192	KDWRGGRTASNIIIPSSTGAAKAVGVKVLPELQGKLTGMAFRVPTDVSVV	241
<a href="#">NP_012542.1</a>	192	KDWRGGRTASNIIIPSSTGAAKAVGVKVLPELQGKLTGMAFRVPTDVSVV	241
<a href="#">XP_456022.1</a>	191	KDWRGGRTASNIIIPSSTGAAKAVGVKVLPELQGKLTGMAFRVPTDVSVV	240
<a href="#">NP_596154.1</a>	194	KDWRGGRGASANIIPSSTGAAKAVGVKIPALNGKLTGMAFRVPTPDVSVV	243
<a href="#">NP_595236.1</a>	194	KDWRGGRGASANIIPSSTGAAKAVGVKIPALNGKLTGMAFRVPTPDVSVV	243
<a href="#">XP_003717853.1</a>	192	KDWRGGRGAAQNIIPSSTGAAKAVGVKIPALNGKLTGMSMRVPTANVSVV	241
<a href="#">XP_956977.1</a>	193	KDWRGGRTAAQNIIPSSTGAAKAVGVKIPDLNGKLTGMAMRVPTANVSVV	242
<a href="#">NP_001060897.1</a>	196	KDWRGGRAASFNIIPSSTGAAKAVGVKLPDLNGKLTGMSFRVPTDVSVV	245
<a href="#">NP_001004949.1</a>	152	KLWRDGRGAGQNIIPASTGAAKAVGVIPELNGKLTGMAFRVPTPNVSVV	201

# Multiple Sequence Alignment (MSA)

MSAs can be challenging to create for more diverse members of a protein family such as the transcription factor Pax6 which is present in hugely divergent species

<a href="#"><u>NP_001595.2</u></a>	17	NGRPLPDSTRQKIVELAHSGARPCDISR---ILQTHADAKVQVLDNQNVS	63
<a href="#"><u>XP_003954413.1</u></a>	11	---LYLIDSRELIAEVGTGWQGDEDALWYIFVLXTHADAKVQVLDNQNVS	57
<a href="#"><u>XP_001084865.2</u></a>	51	NGRPLPDSTRQKIVELAHSGARPCDISR---ILQTHADAKVQVLDNQNVS	97
<a href="#"><u>NP_001035735.1</u></a>	17	NGRPLPDSTRQKIVELAHSGARPCDISR---ILQ-----VS	49
<a href="#"><u>NP_001231127.1</u></a>	17	NGRPLPDSTRQKIVELAHSGARPCDISR---ILQTHADAKVQVLDNENVS	63
<a href="#"><u>NP_037133.1</u></a>	17	NGRPLPDSTRQKIVELAHSGARPCDISR---ILQ-----VS	49
<a href="#"><u>NP_990397.1</u></a>	17	NGRPLPDSTRQKIVELAHSGARPCDISR---ILQTHADAKVQVLDNQNVS	63
<a href="#"><u>NP_571379.1</u></a>	36	NGRPLPDSTRQKIVELAHSGARPCDISR---ILQTHADAKVQVLDNENVS	82
<a href="#"><u>NP_524638.3</u></a>	42	NGRPLPDSTRQKIVELAHSGARPCDISR---ILQ-----VS	74
<a href="#"><u>XP_311087.5</u></a>	33	NGRPLPDSTRQKIVELAHSGARPCDISR---ILQ-----VS	65
<a href="#"><u>NP_001024570.1</u></a>	18	NGRPLPDATRQRIVDLAHKGCRPCDISR---LLQ-----VS	50
<a href="#"><u>NP_001006763.1</u></a>	17	NGRPLPDSTRQKIVELAHSGARPCDISR---ILQ-----VS	49

# Multiple Sequence Alignment (MSA)

Some features of proteins are revealed by MSAs

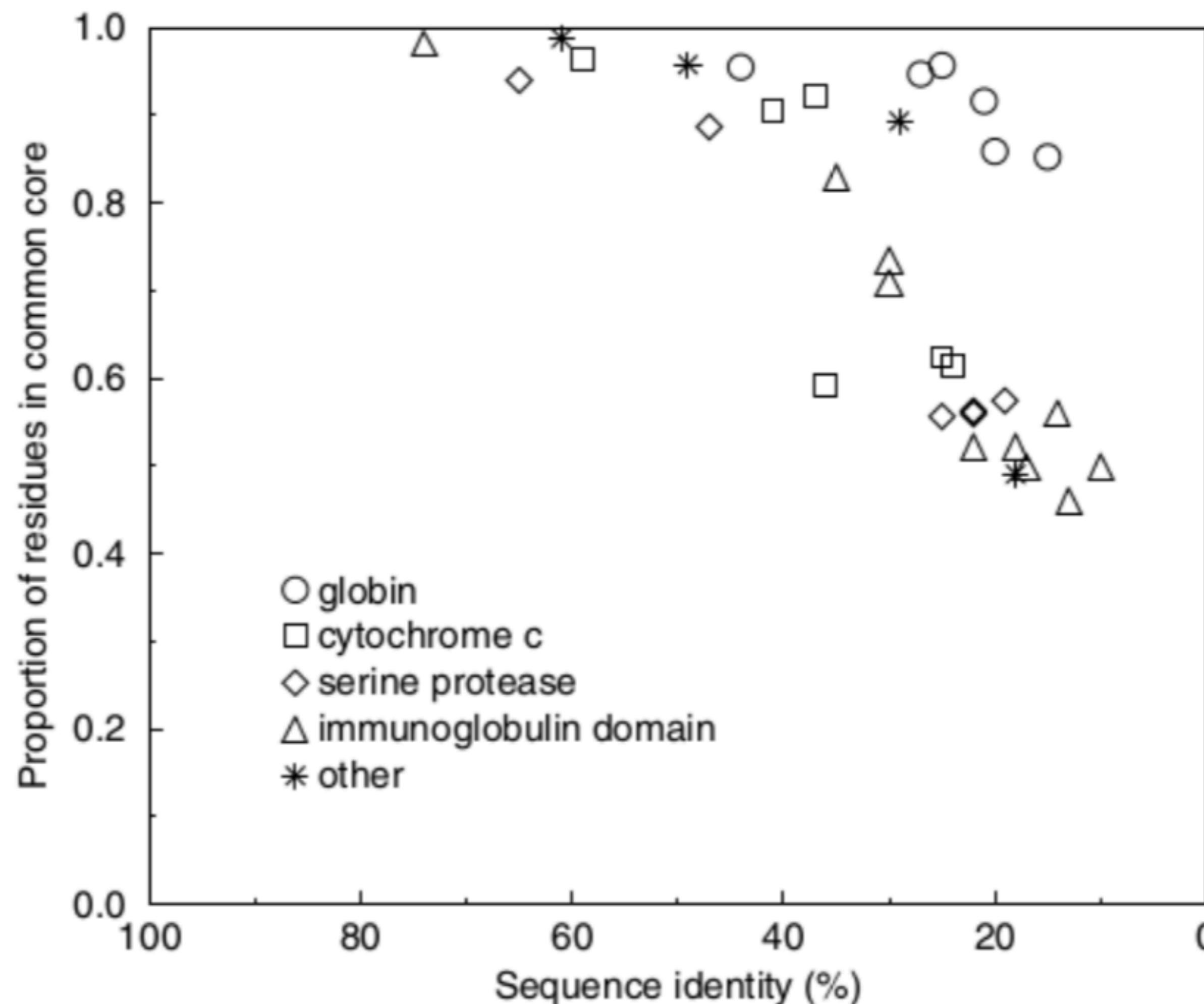
conserved structural features (for example)

- alpha helices
- beta-sheets
- disulfide bridges
- active sites of enzymes
- transmembrane domains

Other parts of sequences may not play a functional/structural role, are not conserved and are not aligned well in MSAs

- linker/spacer regions of proteins
- non-coding, non-regulatory regions of genomic DNA
- dis-ordered parts of proteins

# Sequence vs. Structure



For proteins with significant identity (>30%) where the structure is known most residues are superimposable - there are many different sequences that can converge on similar core structures.

# What can we use MSAs for?

- To infer function. If we have an MSA for a well defined protein family/domain and we can align an unknown sequence to it
- The vast majority of “functional” knowledge about genomes, genes and proteins is actually inferred from such alignments
- As greater numbers of genomes are being sequenced statistical models can be generated such as profile-HMMs and used to describe molecular patterns. These are built from multiple sequence alignments and are crucial to identify proteins present in novel genomes
- These profiles are far more sensitive than pairwise alignment methods such as BLAST. Variants of BLAST (such as DELTA-BLAST) exploit this (NB last week’s lecture!)

# What can we use MSAs for?

- Finding the “functional” parts of sequences
  - such parts are commonly conserved and so align “better” in the MSA than “non functional” parts
- Prediction of potentially deleterious mutation
  - highly conserved residues in MSAs more commonly mark important positions in a regulatory element or protein
- Understanding the evolution of organisms, proteins, genomes
  - MSAs are the first step in performing phylogenetic inference.
- To find possible gene regulatory regions in genomes

# Main Approaches to MSA

- Exact Approaches
- Progressive Sequence Alignment
- Iterative Approaches
- Consistency-Based Approaches
- Structure-Based Approaches

# Exact MSA

In effect performing global pairwise alignment in N-dimensions (NW type Dynamic Programming)

memory complexity  $O_m(L^N)$

time complexity  $O_t(2^NL^N)$

e.g.  $N=10, L=1000: O_m(10^{30}), O_t(10^{33})$

For any reasonably sized set of sequences this is computationally intractable so alternative approaches are needed

# Progressive Sequence Alignment

**Heuristic** approach, not guaranteed to find the optimal alignment

Principle; align sequences pairwise then begin building the MSA starting with the closest pair and then sequentially adding sequences

ClustalW (Thompson et al. 1994) based on progressive alignment idea of Feng & Doolittle (1987,1990)

Step 1 - generate pairwise alignments of all sequences:

	SeqA Name	Len(aa)	SeqB Name	Len(aa)	Score	
1	beta_globin	147	2	myoglobin	154	25
1	beta_globin	147	3	neuroglobin	151	15
1	beta_globin	147	4	soybean	144	13
1	beta_globin	147	5	rice	166	21
2	myoglobin	154	3	neuroglobin	151	16
2	myoglobin	154	4	soybean	144	8
2	myoglobin	154	5	rice	166	12
3	neuroglobin	151	4	soybean	144	17
3	neuroglobin	151	5	rice	166	18
4	soybean	144	5	rice	166	43 ← 1

# Progressive Sequence Alignment

## Scoring the alignments

By default ClustalW just uses percentage identity, but better to calculate distance matrices as these then allow the construction of a “guide tree” to guide the alignment process (i.e. decide in which order to add sequences)

$$D = -\ln S_{eff}$$

$$S_{eff} = \frac{S_{real(i,j)} - S_{rand(i,j)}}{S_{iden(i,j)} - S_{rand(i,j)}} \times 100$$

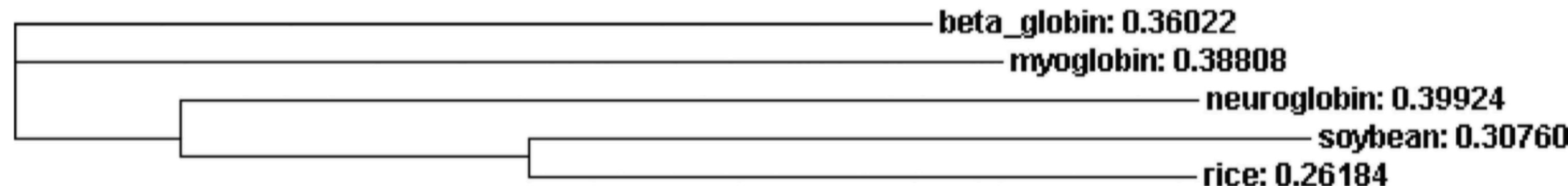
where  $D$  is distance,  $S_{eff}$  is the a normalised similarity score.  $S_{eff}$  is calculated from  $S_{real(i,j)}$  the actual similarity between sequences  $(i,j)$ ,  $S_{iden(i,j)}$ , the average score aligning  $i$  and  $j$  to themselves (in effect a form of length normalisation) and  $S_{rand(i,j)}$  the mean alignment score for many random shuffles of the sequences to be aligned

# Progressive Sequence Alignment

The final distance score depends on  $S_{\text{eff}}$  so that if  $S_{\text{eff}} = 0$ ,  $D = \text{infinite}$  and if  $S_{\text{eff}} = 1$  (they have perfect similarity)  $D = 0$ .

## Step 2 - Guide Tree

Using the Distance Matrix (pairwise distances) create a guide tree for the alignment, for example using the Unweighted Pair Group Method of Arithmetic Averages (UPGMA)



So order :- rice:soybean -> neuroglobin ->beta-globin -> myoglobin

# Progressive Sequence Alignment

## Step 3 - Creation of the MSA

Using the guide tree the closest pair are aligned by DP (e.g. NW) and then the next closest sequence is added to the alignment to produce a profile. This continues sequentially until all sequences have been added

When adding an additional sequence it is essentially a dynamic programming approach, but the gap extension/deletion penalties are dynamically calculated weighted by pre-existing gaps in the more closely related sequence alignments

There are several approaches to this ranging from ones that absolutely preserve gaps in earlier aligned sequences to those that moderate for those so as not to have highly related sequences biasing alignment of more distantly related ones

# Progressive Alignment Example

For example, say your sequences are:

$S_1$	A	T	T	G	C	C	A	T	T
$S_2$	A	T	G	G	C	C	A	T	T
$S_3$	A	T	C	C	A	A	T	T	T
$S_4$	A	T	C	T	T	C	T	T	
$S_5$	A	C	T	G	A	C	C		

Perform the pairwise alignments

$S_1$	A	T	T	G	C	C	A	T	T
$S_2$	A	T	G	G	C	C	A	T	T

$S_1$	A	T	T	G	C	C	A	T	T	-	-
$S_3$	A	T	C	-	C	A	A	T	T	T	T

$S_1$	A	T	T	G	C	C	A	T	T
$S_4$	A	T	C	T	T	C	-	T	T

$S_1$	A	T	T	G	C	C	A	T	T
$S_5$	A	C	T	G	A	C	C	-	-

# Progressive Alignment

Create a Distance matrix and then guide tree (for this simple example just a score matrix)

	S <sub>1</sub>	S <sub>2</sub>	S <sub>3</sub>	S <sub>4</sub>	S <sub>5</sub>	
S <sub>1</sub>	-	7	-2	0	-3	2
S <sub>2</sub>	7	-	-2	0	-4	1
S <sub>3</sub>	-2	-2	-	0	-7	-11
S <sub>4</sub>	0	0	0	-	-3	-3
S <sub>5</sub>	-3	-4	-7	-3	-	-17
	2	1	-11	-3	-17	

S<sub>1</sub> is the sequence most similar to the rest, and below are the best alignments between S<sub>1</sub> and the rest of the sequences.

# Performing the Progressive Alignment

Let's use the alignment of  $S_1$  and  $S_2$ .

$S_1$	A	T	T	G	C	C	A	T	T
$S_2$	A	T	G	G	C	C	A	T	T

$S_1$  and  $S_2$  are aligned

Now, let's add  $S_3$ , using its alignment to  $S_1$ .

$S_1$	A	T	T	G	C	C	A	T	T	-	-
$S_2$	A	T	G	G	C	C	A	T	T	-	-
$S_3$	A	T	C	-	C	A	A	T	T	T	T

$S_1$ ,  $S_2$ , and  $S_3$  are aligned

Then, let's add  $S_4$ , using its alignment to  $S_1$ .

$S_1$	A	T	T	G	C	C	A	T	T	-	-
$S_2$	A	T	G	G	C	C	A	T	T	-	-
$S_3$	A	T	C	-	C	A	A	T	T	T	T
$S_4$	A	T	C	T	T	C	-	T	T	-	-

$S_1$ ,  $S_2$ ,  $S_3$ , and  $S_4$  are aligned

Finally, let's add  $S_5$ , using its alignment to  $S_1$ .

$S_1$	A	T	T	G	C	C	A	T	T	-	-
$S_2$	A	T	G	G	C	C	A	T	T	-	-
$S_3$	A	T	C	-	C	A	A	T	T	T	T
$S_4$	A	T	C	T	T	C	-	T	T	-	-
$S_5$	A	C	T	G	A	C	C	-	-	-	-

$S_1$ ,  $S_2$ ,  $S_3$ ,  $S_4$ , and  $S_5$  are aligned

# Progressive Sequence Alignment

CLUSTAL W (1.83) multiple sequence alignment

beta globin	-----MVHLT <b>PEEKSAVTALW</b> GKVNVD--EVGGEALGRLLVVY <b>PWTQRFFESFG</b> -	47
myoglobin	-----MGLS <b>DGEWQLV</b> LNV <b>WGKVEADIPGHGQEVLIRLFKGHPETLEKFDFK</b> -	48
neuroglobin	-----MERPE <b>PELIRQSWRAVSRS</b> PLEHGTVLFARLF <b>ALEPDLLPLFQYNCR</b>	47
soybean	-----MVAFT <b>EKQDALVSSSFEAFKANIPQYSVVFYTSILEK</b> <b>APAAKDLFSFLA</b> -	49
rice	MALVEDNNNAVAVSFS <b>EEQEALVLKS</b> WAILKKDSANIALRFFLKIFEVAPSASQMFSFLR-	59
	: : : : . . . : * * .	
beta globin	DLST <b>PDAVMGNPKVKAHGKKVLGA</b> FSDGLAHLDNLKGTFATLS---- <b>ELHCDKLHV</b> DPE	102
myoglobin	HLKSEDEM <b>KASEDLKKHGATVLTALGGIL</b> KKKGHHEAEIKPLA----QSHAT <b>KHKIPVK</b>	103
neuroglobin	QFSSPEDCLSS <b>PEFLDHIRKVMLVIDAAVTN</b> VEDLSSL <b>EEYLAS</b> --LGRK <b>HRAVG</b> VKLS	104
soybean	--NGVD <b>PT</b> --NPKLT <b>GHAEKL</b> FALVRDSAGQLKAS <b>GTVVADAA</b> --LGSV <b>HAQKAVTDP</b>	101
rice	--NSDV <b>PLEKNPKLKTHAMSVF</b> VMTC <b>CEAAAQLRK</b> AGKVTVR <b>D</b> TTLKRLGATH <b>HLKY</b> <b>GVGDA</b>	117
	. . . * . : : :	
beta globin	NFRILLGNVLVCVLAHHF-GKEFT <b>PPVQAAYQKV</b> VAGVANALA <b>HKYH</b> -----	147
myoglobin	YLEFISECII <b>QVLIQS</b> KH-PGDFG <b>ADAQGAMN</b> KALELFRKDMASNY <b>KELGFQG</b>	154
neuroglobin	SFSTV <b>GESLLY</b> MLEKCL-GPAFT <b>PATRAAWS</b> QLYGAVV <b>QAMSRGW</b> DGE----	151
soybean	QFVVVK <b>EALLKT</b> IKAAV- <b>GDKWS</b> DELSRAWEVAYDELAA <b>AIKK</b> -----	144
rice	HFEVV <b>KFALLDT</b> IKEE <b>VPADMWS</b> PAMKSAW <b>SEAYDHLVAAIK</b> QEMKPAE---	166
	: : : : : * . . . :	

\* (100% conserved), : (conservative substitution), . (less conservative subst.), and arrowheads (highly conserved residues)

# Iterative Alignment

Principle; create an initial progressive alignment and then optimise it using dynamic programming and a scoring function

Unlike progressive alignment itself, iterative alignment can correct errors made during the initial alignment. Poor guide tree, misplaced gaps etc.

Optimisation needs a scoring function for example the sum-of-pairs score:

M	Q	P	I	L	L	L	V
M	L	R	-	L	L	-	-
M	K	-	I	L	L	L	-
M	P	P	V	L	I	L	V

For each column calculate the pairwise score for every combination. Then sum across all columns.

So for Column 3; P->R, P->-, P->P, R->-, R->P, - ->P

[given a simple scoring scheme of m=+1, mm=-1, gap=-2]

Pair Score = sum(-1, -2, +1, -2, -1, -2) = -6

Summing all pair scores across columns gives -20

# Iterative Alignment

## MUSCLE - 2004

High accuracy and exceptional speed; in the original paper 1000 protein sequences of 282aa were aligned in 21 seconds on a standard PC

Step 1 - Generate a draft progressive alignment, this can use identity or k-mer profiling of sequences; e.g. frequency fingerprints for k-mers between sequences (very fast)

ACGCTACGCATACGCATG

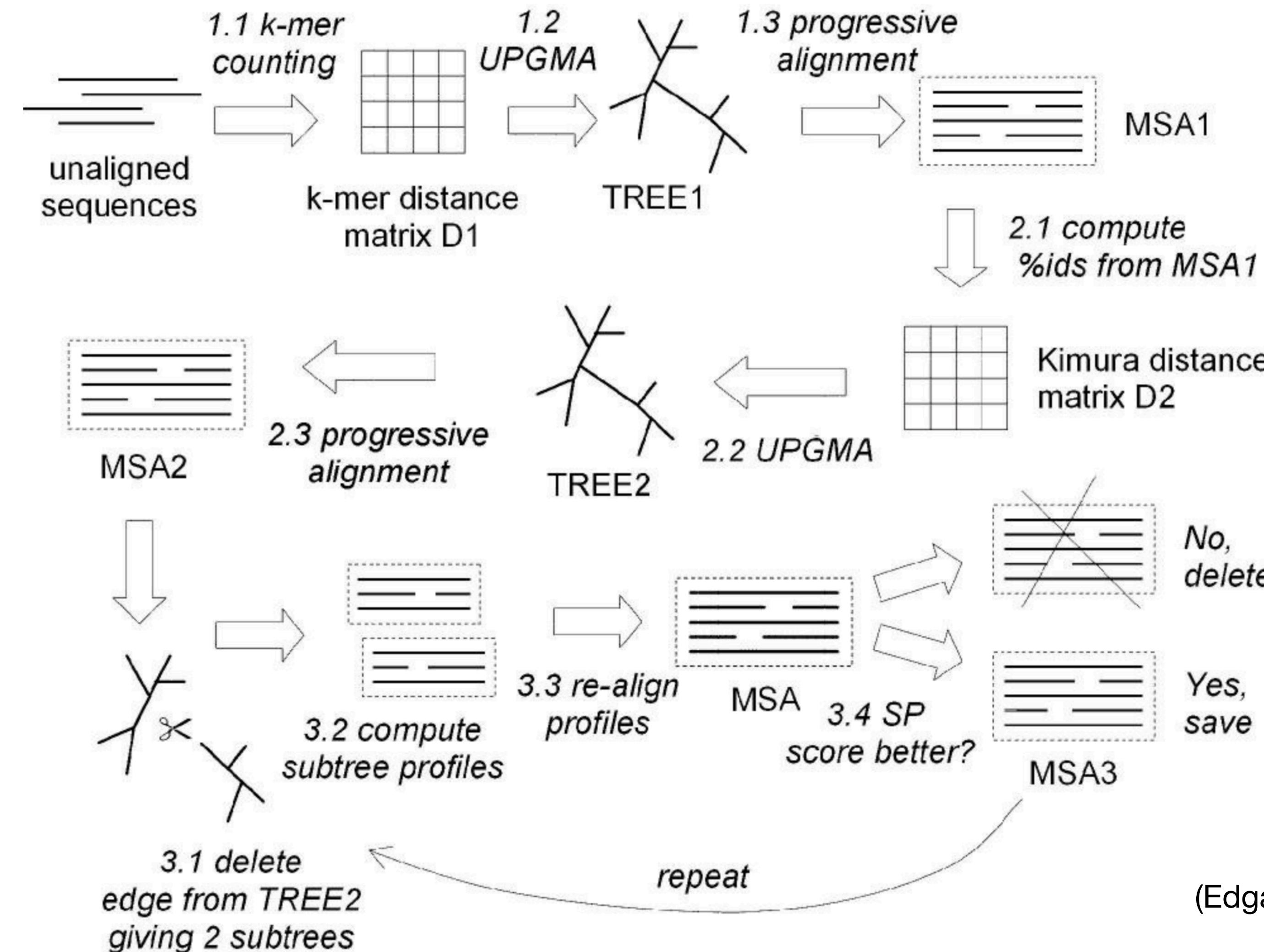
ACG TAC CAT CGC  
CGC ACG ATA GCA  
GCT CGC TAC CAT  
CTA GCA ACG ATG

k-mer profiling (3-mers)

A guide tree is then built from either of these using for example UPGMA. Sequences are progressively added to the alignment as before

# Iterative Alignment

Schematic of **MUSCLE** methodology



# Iterative Alignment

Multiple Sequence Comparison by Log Expectation (**MUSCLE**) (Edgar 2004)

Step 2 - MUSCLE next improves the tree and builds a new progressive alignment. It does this by calculating a new distance matrix by using the Kimura distance (this attempts to model the possibility of multiple changes in sequence happening at the same position).

Using the new distance matrix a new tree is created and a new progressive alignment made.

Step 3 - Re-partitioning. Next for every branch in the tree MUSCLE creates sub-partitions and creates profiles (sub-alignments) of each and these two profiles are then aligned to create a new MSA. Consider this a set of local MSAs testing every possible bi-partition of the tree.

For every tree created this way the PSP is compared to the previous best. If it is worse the MSA from that bi-partition is deleted if better it replaces the MSA.

$$d_{\text{Kimura}} = -\log_e (1 - D - D^2/5)$$

# MUSCLE profile-profile Alignment

In practice sum of pair (SP) scores are normally calculated as probabilistic scores and scoring functions use scoring matrices such as PAM, BLOSUM.

$$S_{ij} = \log(p_{ij}/p_i p_j)$$

$$PSP^{xy} = \sum_i \sum_j f_i^x f_j^y S_{ij}$$

For amino acid types i and j,  $p_i$  is the background probability of i,  $p_{ij}$  is the joint probability of i and j being aligned,  $S_{ij}$  is the score from a substitution matrix,  $f_i^x$  is the observed frequency of i in column x of the first profile, PSP is a sequence-weighted sum of substitution matrix scores for each pair of letters (one from each column that is being aligned in a pairwise fashion). The PSP function maximises the sum-of-pairs objective score.

# Iterative Alignment

**MUSCLE** uses a log expectation (LE) score that is in effect a PSP modified to score profile-to-profile alignments

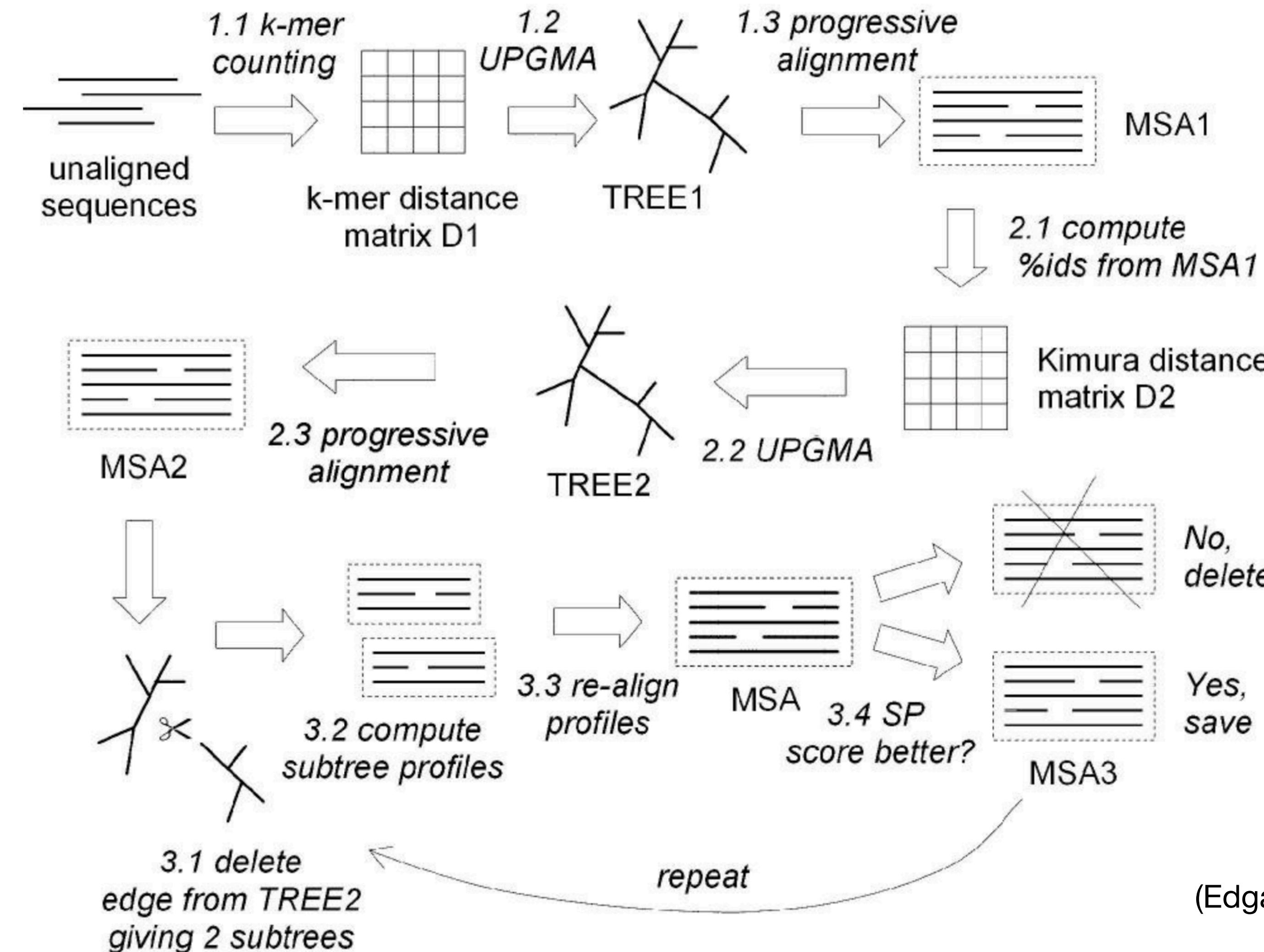
$$LE^{xy} = (1 - f_G^x)(1 - f_G^y) \log \sum_i \sum_j f_i^x f_j^y \frac{p_{ij}}{p_i p_j}$$

For profiles X and Y the LE is calculated. Columns are compared between profiles.  $(1 - f_G^x)$  is the occupancy of a column of profile X. Note the summed term from PSP is logged. This has been shown to increase accuracy of the alignment by rewarding columns that are full (better aligned) over those that have gaps.

**MUSCLE** uses PAM matrices in its PSP function to better refine the amino acid pair scoring.

# Iterative Alignment

Schematic of **MUSCLE** methodology



# Progressive Sequence Alignment

CLUSTAL W (1.83) multiple sequence alignment

beta globin	-----MVHLT <b>PEEKSAVTALW</b> GKVNVD--EVGGEALGRLLVVY	PWTQRFFESFG-	47
myoglobin	-----MGLS <b>DGEWQLVLN</b> VGKVEAD <b>I</b> PGHGQEVLIRLFKGH	PETLEKFDKFK-	48
neuroglobin	-----MERPE <b>PELIRQSWRAVSRS</b> PLEHGTVLFARLFALEPDLLPLFQYNCR		47
soybean	-----MVAFT <b>EKQDALVSSSFEAFKANIPQYSVVFYTSILEKA</b> PAAKDLFSFL <b>A</b> -		49
rice	MALVEDNNNAVAVSFS <b>EEQEALVILKSWAILKKDSANIALRFFLKIFEVAPSASQMFSFLR</b> -		59

: : : : . . . . : \* \* .

beta globin	DLST <b>PDAVMGNPKVKAHGKKVLGAFSDG</b> LAHLDNLKGTF <b>ATLS</b> -----	<b>ELHCDKLHVDPE</b>	102
myoglobin	HLKSEDEM <b>KASEDLKKHGATVLTALGGIL</b> KKKGHHEAEIKPLA-----	<b>QSHATKHKIPVK</b>	103
neuroglobin	QFSSPEDCLSS <b>PEFLDHIRKVMLVIDAAVTN</b> VEDLSSLEEYLAS---	<b>LGRKHRAVGVKLS</b>	104
soybean	-- <b>NGVDPT</b> --NPKLTGHAEKLFALVRDSAGQLKAS <b>GTVVADAA</b> ---	<b>LGSVHAQKAVTDP</b>	101
rice	--NSDVPLEKN <b>PKLKTHAMSVFVMTCEAAAQLRKAGKVTVRDTTLKRLGATHLKY</b> GVGDA		117

. . . \* . : : : :

beta globin	NFRLLGNVLVCVLAHHF-GKEFT <b>PPVQAAYQKV</b> VAGVANALAHKYH-----	147
myoglobin	YLEFISECIIQVLQSKH-PGDFG <b>ADAQGAMNKALELFRKDMASNYKEL</b> GFQG	154
neuroglobin	SFSTVGESLLYMLEKCL-GPAFT <b>PATRAAWSQLYGA</b> VVQAMSRGWDGE---	151
soybean	QFVVVKEALLKTIKAAV- <b>GDKWSDELSRAWEVAYDELAAAIKK</b> A-----	144
rice	HFEVVKFALLDTIKEEVPADMWS <b>PAMKSAWSEAYDHLVAAIKQEMKPAE</b> ---	166

: : : : : \* . . . :

\* (100% conserved) : (conservative substitution) . (less conservative subst.), arrowheads (highly conserved residues)

# Iterative Alignment

MUSCLE (3.6) multiple sequence alignment

beta globin	-----MVHLT <b>PEEKSAVTALW</b> GKVNVD--EVGGEALGRLLVVY	PWTQRFFES-FG
myoglobin	-----MGLS <b>DGEWQLVLNVW</b> GKVEADIPGHGQEVLIRLFKGH	PETLEKFDK-FK
neuroglobin	-----MERPE <b>PELIROQSW</b> RAVSRS	PLEHGTVLFARLFALEPDLLPLFQYNCR
soybean	-----MVAFT <b>EKQDALVSSSFEAFKAN</b> IPQYSVVFYTSILEKAPA	<b>AKDLFSF-LA</b>
rice	MALVEDNNNAVAVSF <b>EEQEALVLKSWAILKKD</b> SANIALRFFLKIFEVAPSASQMFSF-LR	
	: : : : .. .	:: * *
beta globin	DLST <b>PDAVMGNPKVKAHGKKVLGAF</b> --SDG	LAHLDNLKGTF <b>ATLSELHCDKLH</b> --VDPE
myoglobin	HLKSEDEM <b>KASEDLKKHGATVLTAL</b> --GGILKKKGHHEAEI	<b>KPLAQSHATKHK</b> --IPVK
neuroglobin	QFSSPEDCLSS <b>PEFLDHIRKVMLVI</b> --DAAVTNVEDLSSLEEYLASLGRKHRAV	GVKLS
soybean	<b>NGVDP</b> --TN <b>PKLTGHAEKL</b> FALVRDSAGQLKASGTVVAD--AALGSVH <b>AQKAVTDP</b>	
rice	NSDVP--LEKN <b>PKLKTHAMSVFVMTCEAAAQLRK</b> AGKVTVRDTTLKRLGATHLYGVGDA	
	. . . * . : :	: :
beta globin	<b>NFRLLGNVLVCVLAHHGKE-FTPPVQAAYQKV</b> VAGVANALA <b>HKYH</b> -----	
myoglobin	<b>YLEFISECIIIQVLQSKH</b> PGD-FGADAQGAMNKALELFRKDMASNY <b>KELGFQG</b>	
neuroglobin	SFSTVGESLLYMLEKCLGPA-FT <b>PATRAAWSQLY</b> GAVVQAMSRGWDGE-----	
soybean	<b>QFVVVK</b> KEALLKTIKAAVGDK-WS <b>DELSRAWEVAY</b> DELAAAIKKA-----	
rice	<b>HFEVVKFALLDTIKEE</b> VPADMWS <b>PAMKSAWSEAY</b> DHLVAAIK <b>QEMKPAE</b> ---	
	: : : : : * . . . :	

# Consistency-Based Approaches

Principle; if for sequences  $x$ ,  $y$  and  $z$  if amino acid  $x_i$  aligns with  $z_k$  and  $y_j$  aligns with  $z_k$  then  $x_i$  should align with  $y_j$

This can be expressed as a conditional probability:-

$$P(x_i \sim y_j | x, y)$$

$$P(x_i \sim y_j | x, y, z) \approx \sum_k P(x_i \sim z_k | x, z) P(y_j \sim z_k | y, z)$$

It means that MSAs are adjusted as sequences are being added based on the "column context" of the alignment. This is a fundamentally different approach than we have seen so far.

# Consistency-Based Approaches

## Probabilistic Consistency-Based Alignment (**ProbCons**)

Step 1 - ProbCons first uses a paired-hidden Markov Model to calculate the posterior probability matrices for every pair of sequences to make up the MSA (i.e. a probabilistic scoring matrix for each alignment)

Step 2 - Using a NW type approach pairwise alignments are made using the posteriors from Step1, these are then refined using consistency information harvested from across all the pairwise alignments

Step 3 - A guide tree is then made and progressive alignment performed in much the same way as in previous methods.

Step 4 - This process can be iterated until the MSA scores have converged.

# Consistency-Based Alignment

PROBCONS

beta globin	M-----VHLT <b>PEEKSAVTALW</b> GKVNVD--EVGGEALGRLLVVY	PWTQRFFES-FG
myoglobin	M-----GLS <b>DGEWQLVLNVW</b> GKVEAD	I PGHGQEVLIRLFKGHPETLEKFDK-FK
neuroglobin	M-----ERPE <b>PELIRQSW</b> RAVSRS	PLEHGTVLFARLFAL
soybean	M-----VAFT <b>EKQDALVSSS</b> FEAFKAN	I PQYSVVFYTSILEKAP <b>A</b> KDLFSF-LA
rice	MALVEDNNAAVASF <b>EEQEALVLKS</b> WAILK	KD <b>SANIALRFFLKIFEV</b> APSASQMFSF-LR
	*	*
beta globin	DLST <b>PDAVMGNPKVK</b> AHGKKVLGAFSDG	LAHLD---NLK---GTF <b>ATLSEL</b> HCDKLHVDP
myoglobin	HLKSEDEM <b>KAS</b> EDLKKHGATVLTALGGI	---LKKKGHE---AEI <b>KPLAQSHAT</b> KHKIPV
neuroglobin	QFSSPEDCLSS <b>PEFLDH</b> IRKVMLVIDAAVTNVEDLSSLE	---EYLASLGRKHRAV-GVKL
soybean	<b>NGVDP</b> ---TN <b>PKLTGHAEKL</b> FALVRDSAGQLKASGTVV	---AD <b>AALGSVH</b> AQK-AVTD
rice	NSDVP--LEKN <b>PKLKTHAMSVF</b> VMTCEAAAQLRK	AGKVTVRDTTLKRLGATHLY-GVGD
	.	*
beta globin	<b>ENFRLLGNVLVCVLAHHF</b> -GKEFTPPVQAAYQKV	VAGVANALAHK-----YH
myoglobin	<b>KYLEFISECIIQVLQSKH</b> -PGDFGADAQGAMNKALELFRKDMASNYKEL	GFQG
neuroglobin	SSFSTVGESLLYMLEKCL-GPAFT <b>PATRAAWSQLY</b> GAVVQAMSRG	---W-DGE
soybean	<b>PQFVVVK</b> EALLKTIKAAV-GDKWSDELSRAWEVAYDELAAAIK	-----KA
rice	AHFEVVKFALLDTIKEEVPADMWS <b>PAMKSAWSEAY</b> DHLVAAIKQE	---MKPAE
	:	:

# Structural-Based Approaches

Principle; incorporation of protein structural information can improve MSAs

Sequences to be aligned are first used in BLAST searches of the Protein Data Bank (PDB). Matches are then used to guide the creation of the multiple sequence alignment.

The Espresso tool from the T-COFFEE alignment suite does this and works on the principle that tertiary structure of proteins evolve much more slowly than lower order features.

Ultimately this allows for different amino acids producing similar higher order protein structures (which can commonly happen). These structural approaches try to capture that by incorporating actual crystal structure data.

Whilst powerful, this approach is limited by the available crystal structures for proteins.  
[NB new structure prediction tools such as AlphaFold2 are changing this]

# Databases of MSAs

A lot of databases exist where MSAs for protein families have already been computed using some of the tools we have looked at:-

PFAM

<http://pfam.xfam.org/>

Protein domains/families identified by profile-HMM

SMART - Simple Modular Architecture Research Tool

<http://smart.embl-heidelberg.de>

CDD – Conserved Domain Database

<https://www.ncbi.nlm.nih.gov/cdd>

# PFAM protein family MSA

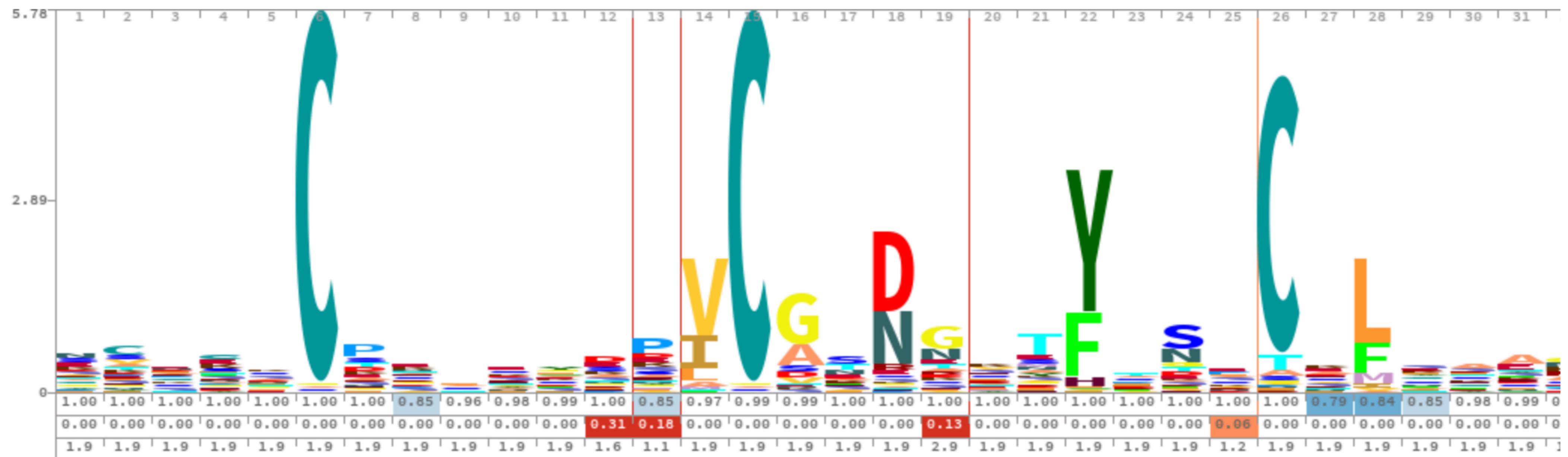
## Kazal2 - Serine Protease Inhibitor

IACS\_PIG/43-88  
CO6\_HUMAN/788-837  
Q95011\_CAEEL/325-376  
HTRA4\_HUMAN/103-152  
CFAI\_RAT/65-109  
Q9VUD9\_DROME/293-342  
Q9VUD8\_DROME/255-310  
AGRIN\_RAT/92-137  
HTRA3\_MOUSE/82-132  
Q95011\_CAEEL/133-183  
FSTL3\_MOUSE/118-165  
FST\_PIG/117-164  
THBI\_RHOPR/55-101  
ISK5\_HUMAN/106-151  
IOV7\_CHICK/302-347  
F1PR78\_CANLF/19-64  
IOVO\_MELGA/81-126  
Q9W271\_DROME/465-512  
Q9W269\_DROME/541-588  
SO2A1\_MOUSE/444-493  
SO1B3\_HUMAN/460-506  
SO1B2\_MOUSE/455-503  
SO1C1\_HUMAN/477-523  
SO1A1\_RAT/440-486  
SO1A2\_HUMAN/440-486  
SO4A1\_HUMAN/505-553  
Q8SY02\_DROME/481-527  
Q9D5W6\_MOUSE/481-527  
A0A0G2K9W9\_RAT/488-534  
SO5A1\_HUMAN/556-603  
Q9VP16\_DROME/514-563  
TICN3\_HUMAN/137-183  
TICN2\_MOUSE/134-180  
TICN1\_HUMAN/134-180  
FST\_PIG/191-239  
FSTA\_DANRE/272-319  
FSTL1\_MOUSE/51-96  
SPRL1\_MOUSE/443-495  
SPRC\_RABIT/98-151  
SPRC\_CAEEL/75-135

SHLFFCT.REMD..P.ICGTNG...KSYANP.CIFCSEKLRNEKFD.....FGHWGHC  
SPEEDCSHHSED....LCVFDTDSNDYFTSPACKFLAEKCLNNQQQLH.....FLHIGSC  
STCITCPKDEKK.IP.ICDNRN...MTHPTL.CSFIQYNCEARNNED..EERV.....LVHIKSC  
PSTCGCP TLGGA....VCGSDR...RTYPSM.CALRAENRAARRLGKVPAV.....PVQWGNC  
KLPYQCP..KAG.TP.VCATNG...RGYPTY.CHLKSFEC.LHPEIK.....FSNNGTC  
KCNFQCPKASLS....ICASNGKCVVNFPGQ.CELSQWNCFNTKNVFH.....QVHDAEC  
NCDELCKFEYS..P.ICAHNGICIHEFANQ.CVMNTFNCKHRDL SFRAVDED.....VCRLGVC  
CKKNACPATVAP....VCGSDA...STYSNE.CELQRAQCNCQRRIR.....LLRQGPC  
CVRGVCR.CRWT..HTVCGTDG...HTYADV.CAL.QAASRRALQVSGTPV R.....QLQKGAC  
DCNHNCNTTEFD..P.VCDTNG...SVYRNL.CVFQMRRC ELOLESQR..IQL.....AEDRKFC  
ECVPNCEGLPAG.FQ.VCGSDG...ATYRDE.CELRTARCRGHPDLR.....VMYRGRC  
VCAPDCSNITWK.GP.VCGLDG...KTYRNE.CALLKARCKEQPELE.....VQYQGKC  
DVCQECDGDEYK..P.VCGSDD...ITYDNN.CRLECASISSSPGVE.....LKHEGPC  
DGDFICP.DYYE..A.VCGTDG...KTYDNR.CALCAENAKTGSQIG.....VKSEGE C  
EAITACP FILQE....VCGTDG...VTYSND.CSLCAHNIELGTSVA.....KKHDGRC  
GSQIACP.RHLQ..P. ICGTDH...KTYSNE.CMFCA TLNKKFEVR.....KLQDTAC  
KVMILCN.KALN..P.VCGTDG...VTYDNE.CVLCAHNLEQGTSVG.....KKHDGEC  
SASCHCDYVHYA..P.VCSADN...ITFISA.C...HAGCSERTKDA..LGRT.....IYTGCEC  
NSACSCDYVRY S..P.VCGENN...MTYISA.C...HAGCKKLLVNSE.GKK.....IFYDCSC  
RRDCLCPDSVFH..P.VCGDNG...VEYLSP.C...HAGCSSLNVSSAASKQP.....IYLN CSC  
NSENCNCDESQWE..P.VCGNNG...ITYLSP.C...LAGCKSSSGIK...KHT.....VFYNCSC  
NSDCICDKNQWE..P.VCGENG...VTYISP.C...LAGCKSFRGDKKLMNI.....EFYDCSC  
NSRCKCSETKWE..P.MCGENG...ITYVSA.C...LAGCQTSNRSG...KNI.....IFYNCTC  
NTRCSCSTNTWD..P.VCGDNG...VAYMSA.C...LAGCKKFVGTG...TNM.....VFQDCSC  
NVDCNCPSKIWD..P.VCGNNG...LSYLSA.C...LAGCETSIGTG...INM.....VFQNCSC  
NAACSCQPEHYS..P.VCGSDG...L MYFSL.C...HAGCPAATE TNVDGQK.....VYRDCSC  
NSNC GCSR TNYD..P. ICGVDG...VMYYSP.C...YAGCVQEEHAN...SLK.....RYHNCSC  
NPVPGCTTSEYN..P.VCGRDE...TQYFSP.C...FAGCKATKKLR...KEK.....TYYNCSC  
NYHCACTTSLYS..S.VCGRDE...KEYFSP.C...FAGCSATKVQQ...NEK.....TYYNCSC  
NVNC GCKI HEYE..P.VCGSDG...ITYFNP.C...LAGCVNSGNLST.GIR.....NYTECTC  
YCEKICANVYDENDEII CGSDG...YMYTGE.TQLQCYSSCLNISVT.....IKSKGSC  
STCKQCPVVYPS..P.VCGSDG...HTYSFQ.CKLEYQACVLGKQIS.....VKCEGH C  
SVCKPCHMAQLA..S.VCGSDG...HTYSSV.CKLEQQACLSSKQLA.....VRCEGPC  
VKCKPCPVAQSA..M.VCGSDG...HSYTSK.CKLEFHACSTGKSLA.....TLCDGPC  
TCNRICPEPTSSEQY.LCGNDG...VTYSSA.CHLRKATCLLGRSIG.....LAYEGKC  
VCAESCPESRSE.EA.VCASDN...TTYPSE.CAMKQAACSLGVILLE.....VKHSGSC  
LCIEQCK.PHKR..P.VCGSNG...KTYLNH.CELHRDACL TGSKI Q.....VDYDGHC  
QDPETCPPAKIL.DQ.AC GTDN...QTYASS.CHLFATKCRLEGTKKGHQLQ.....LDYFGAC  
QDPTSCPAPVGEFEK.VCSNDN...KTFDSS.CHFFATKCTLEGTKKGHKLH.....LDYIGPC  
ECISKCP ELDGDPMDKVCANNN...QTFTSL.CDLYRERCLCKRSKECSKA FNAKVHLEYLGEC

# PFAM family profile-HMM

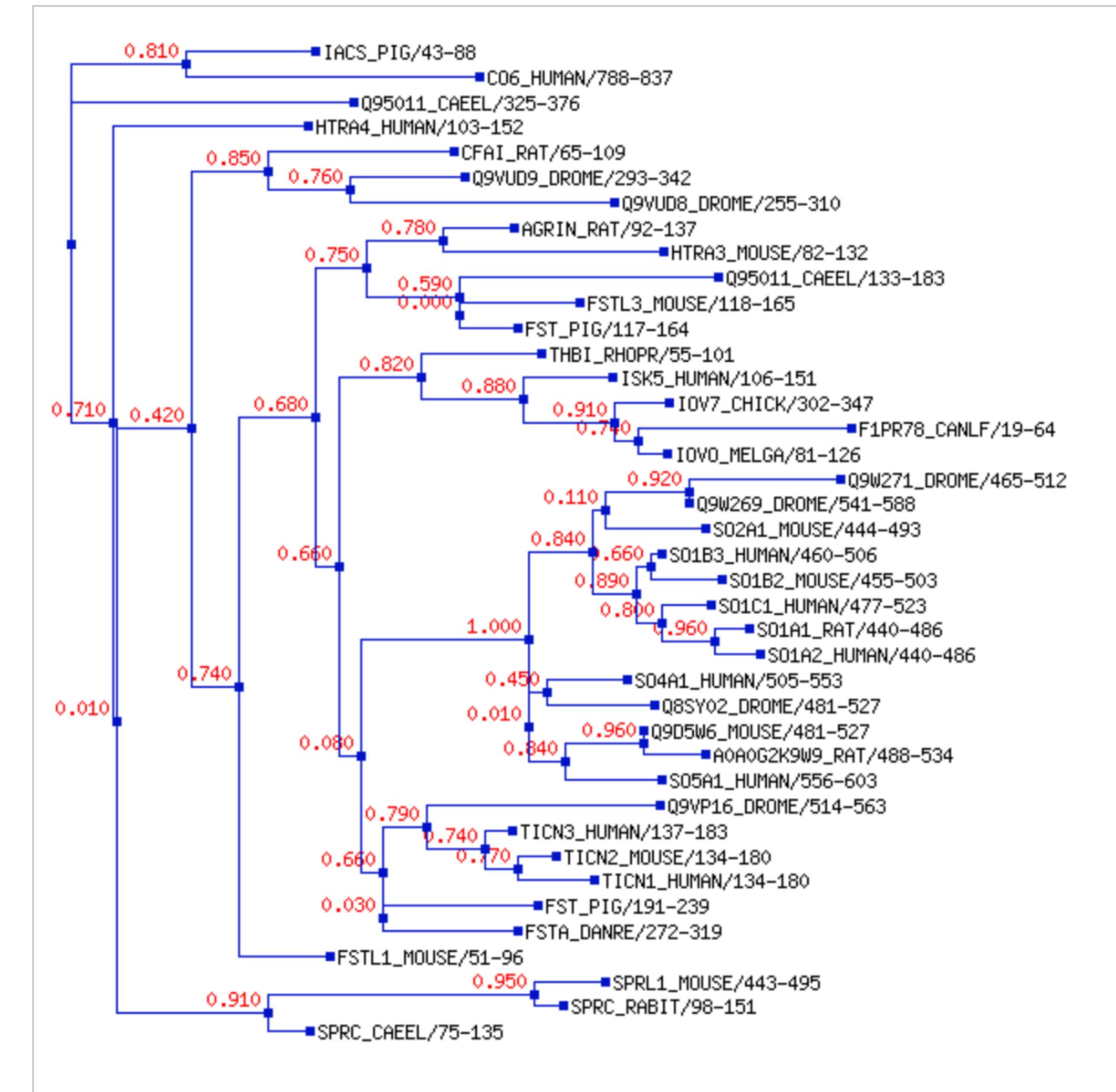
## Kazal2 - Serine Protease Inhibitor



sequence logo (plotting information content per column)

# PFAM Tree

## Kazal2 - Serine Protease Inhibitor



# Further Reading (optional)

## MSA

- Feng, D. F., and Doolittle, R. F. Progressive sequence alignment as a prerequisite to correct phylogenetic trees. *J. Mol. Evol.* 25, 351–360 (1987).
- Doolittle, R. F. On the trail of protein sequences. *Bioinformatics* 16, 24–33 (2000).
- Higgins, D. G., Blackshields, G., and Wallace, I. M. Mind the gaps: Progress in progressive alignment. *Proc. Natl. Acad. Sci. USA.* 102, 10411–10412 (2005).
- Kumar, S., and Filipski, A. Multiple sequence alignment: in pur- suit of homologous DNA positions. *Genome Res.* 17, 127 – 135 (2007).

## CLUSTALW

- Thompson, J. D., Higgins, D. G., and Gibson, T. J. CLUSTALW: Improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res.* 22, 4673–4680 (1994).

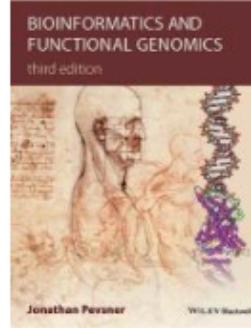
## MUSCLE

- Edgar, R. C. MUSCLE: A multiple sequence alignment method with reduced time and space complexity. *BMC Bioinformatics* 5, 113 (2004a).
- Edgar, R. C. MUSCLE: Multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.* 32, 1792–1797 (2004b).

## PFAM

- Bateman, A., Coin, L., Durbin, R., Finn, R. D., Hollich, V., Griffiths-Jones, S., Khanna, A., Marshall, M., Moxon, S., Sonnhammer, E. L., Studholme, D. J., Yeats, C., and Eddy, S. R. The Pfam protein families database. *Nucleic Acids Res.* 32(Database issue), D138–D141 (2004).

## Week 5 - Exploring Biological Databases



This week you should browse BFG for examples of some of these databases in use, especially material about biological databases in Chapter 2, but you don't need to read the whole chapter.

Other useful things you might like to browse are some of the useful guides provided by some of the resources above such as:-

- NCBI Training Tutorials - <https://www.ncbi.nlm.nih.gov/guide/training-tutorials/> (very good)
- PubMed User Guide - <https://pubmed.ncbi.nlm.nih.gov/help/> (very good)
- Bioportal Help - [https://www.bioontology.org/wiki/BioPortal\\_Help](https://www.bioontology.org/wiki/BioPortal_Help)
- Biogrid Help - <https://wiki.thebiogrid.org>
- Reactome User Guide - <https://reactome.org/userguide> (very good)