

### Multiple alignment



## Orthologous genes vs. Parologous Genes

Orthologous genes: are homologous genes, present in different but related species, which encode for proteins that have similar functions and that have been separated not by a duplication event, but as a result of speciation (separation of species). (Homologous genes, different in different organisms derived from specialization of the species, encode the same function)

Example:  $\alpha$ -globin proteins from human and mouse have diverged about 80 million years ago, when there was the division that gave rise to primates and rodents. The two genes are considered orthologous.

<u>Paralogous genes</u>: are genes originated from the duplication of a single gene in the same organism. (Different genes generated by gene duplication of an organism)

es.  $\alpha$ -globin and  $\beta$ -globin human started to diverge after the duplication of an ancestral globin gene. The two genes are considered paralogs.



#### **Multiple alignment**

The best way to know the characteristics of a particular proteins family is to align many proteins with similar function.

The functionally or structurally most relevant sites tend to remain unchanged in homologous proteins, while the less important sites may also change a lot.

Observe and study the <u>conservations</u> means understanding <u>how</u> <u>the protein families work</u>, what makes them different from each other, <u>whether or not</u> there exist inter- and intra-family <u>phylogenetic relationships</u>.

In this way it is possible to **identify the function of an unknown protein only by observing the sequence of its residues**.

The sequences to multi-align typically are obtained from the database search using the search systems for similarities such as BLAST and FASTA.

Since they come from an alignment (even if produced with heuristic methods) and that they take into account only sequences that have a high score (and a low E, expectation value), the alignment of these multiple DATASET will give satisfactory results.

In a multiple alignment there are taken into account the columns of residues, rather than the proteins to which they belong. Any residue is lined up to be implicitly considered as evolutionarily related, in some way.



# Biological Meaning of the multiple alignment

The multiple alignment summarizes:

- The <u>evolutionary</u> history of a family of proteins.
- The conservation of <u>function</u> related residues.
- The conservation of the <u>structure</u> related residues.

# Methods for the multiple alignment

optimal alignment

-dynamic programming, MSA

- heuristic alignment
  - -Progressive
    - •Global (CLUSTALW, Pileup)
    - •Local (PIMA)
  - Iterative
    - •Global (PRRP)
    - Local (dialign)



# Optimal Multi-alignment Complexity

On the surface of the cube you will have the comparison matrices in pairs between the sequences A-B, B-C and A-C. The optimum alignment of three sequences (A-B-C) requires the filling of the cube and the evaluation of all possible moves within the cube.



The complexity of this algorithm is  $(O(L^N))$ , where L is the length and N is the number of sequences. For three sequences of 300 amino acids the number of comparisons is 2.7 x 10<sup>7</sup>. A complete algorithm of dynamic programming is applicable only in the case of three sequences.

# Optimal Alignment with dynamic programming

Carrillo & Lipmann, 1988



To find an optimal alignment between the three sequences is necessary to calculate the scores within the volume in gray: the volume is bounded by the projections of the areas on the faces of the cube.

This algorithm (implemented in the MSA program) can be used for a few (<10) short sequences.



## Heuristic Algorithm: Progressive Alignment

Alignment of the three sequences A, B, C for the subsequent steps









## **Progressive Alignment**

- Alignments in pairs [N (N-1) / 2 comparisons] of all sequences with dynamic programming or approximate methods (BLAST, FASTA). Calculate a diagonal matrix of distances.
- Build a tree (Neighbor-joining, UPGMA, etc.) on the basis of the matrix. The tree serves as a guide for subsequent alignments.

• Beginning with the first node added to the tree, which represents the closest two sequences, align gradually the other nodes (which may be two sequences, a sequence and an alignment or two alignments) until all sequences have been aligned

# **Distance Matrix Calculation**

Objective: From the single alignment performed at the first step, we have the score representing the similarity between the aligned pairs. From it we compute the edit distance for each pair.

- now very important: scoring
- so far: type 1: similarity (e.g., PAM)
  - type 2: distances  $\hat{=}$  metric (edit distance) *minimum entropy*

• however: distances and similarities can be translated into each other.

- often: given similarity (PAM, BLOSUM), searched distance
- one of the first approaches: Feng-Doolittle
  - given: 2 sequences a,b and a similarity function  $S(\cdot, \cdot)$
  - wanted: 2 sequences a, b and a distance function  $D(\cdot, \cdot)$



# **Distance Matrix Calculation**

#### $\label{eq:Feng-Doolittle: Similarity} \leftrightarrow \mbox{Distance}$

#### • definitions:

• S<sub>rand</sub>

• S(a, b)•  $S_{a,b}^{\max} = \frac{S(a, a) + S(b, b)}{2}$  similarity of *a*, *b* 

maximal possible similarity

expected score for aligning two random sequences a', b' with same length and compositions  $\Rightarrow a, b$  shuffled

e,

- then: define  $S_{a,b}^{\text{eff}} = \frac{S(a,b) S_{rand}}{S_{a,b}^{max} S_{rand}}$
- $S_{a,b}^{\text{eff}}$  normalised percentage similarity
  - values between 0 and 1
  - converges exponentially slow to 0 for increasing evolutionary distance
- hence: for an approx. linear distance define  $D(a, b) = -\log S_{a,b}^{eff}$
- comment: D(a, b) will be later used to build guide tree in "progressive MSA"

## The guide Tree for the clusterization

- Give a score for each pair, it is possible to determine which is the best way to pass from a protein to an another following a specific order that can be considered as evolutionary progression.
- The (rapid) methods to build the tree, given a score set, are two:
  - Neighbour Joining (NJ)
  - Unweighted Pair Group Methos with Arithmetics Mean (UPGMA)
- These methods base their topological schema on the Distance Matrices that have different meaning with respect to the scores:
  - **Score**: high value => more similar sequences
  - Distance : high value => more distant sequences => less similar sequences



#### **Clusterization and guide tree**

1 Hbb\_human
2 Hbb\_horse
3 Hba\_human
4 Hba\_horse
5 Myg whale



It is a matrix of distances, the lower the number, the greater the similarity



Guide Tree given by the clusterization order

PEEKSAVTALWGKVN--VDEVGG Hbb\_human GEEKAAVLALWDKVN--EEEVGG Hbb\_horse PADKTNVKAAWGKVGAHAGEYGA Hba\_human AADKTNVKAAWSKVGGHAGEYGA Hba\_horse EHEWQLVLHVWAKVEAGVAGHGQ Myg\_whale

**Final Alignment** 



# CLUSTAL: example

Higgins & Sharp 1988



#### 

Seq1A A A A ASeq2A A A A ASeq3A A A CSeq4A A C C

N(N-1)/2 comparisons



It is used the method of the sum of the score pairs in each column to determine the total alignment score.

This method does not take into account the history of the sequences and the fact that a same character in the column can be easily shared by very similar sequences for evolutionary reasons.



To remove the bias introduced by proteins in the same family, it is assigned a weight to the sequences in order to increase the score in the comparisons of evolutionarily distant sequences and decrease it in comparisons between neighboring sequences. Generally it is added up all the scores of all possible pairs of aligned proteins, weighing the values based on the similarity in the same cluster to prevent some cluster prevail over others in the final count. The WSP (Weighted Sum of Pairs) is defined as:

$$WSP_{score} = \sum_{i=1}^{N-1} \sum_{j=1}^{N} W_{ij} \hat{S}(A_{ij})$$

Objective Function (OF)

N: sequence number i,j: sequence pair \$: similarity score of the pair W: weight of the pair

The total value of the WSP depends on the scoring criteria used in the alignment rather than by biological considerations, but it is a valid criterion for all the alignments with the same parameters

# CLUSTALW improvement



Thompson et al 1994

The most phylogenetically distant sequences receive a weight proportionally higher in the alignment

FIG. 3. Sequence weights for the seven globin sequences from Fig. 1. A rooted neighborjoining tree is shown with branch lengths. The weights are shown for each sequence before normalization (the weights are normalized so as to make the largest equal to 1.0).

#### **RESIDUE-SPECIFIC GAP UPENING PENALTY FACTORS"**

Residue	Penalty	Residue	Penalty
A	1.13	М	1.29
С	1.13	Ν	0.63
D	0.96	Р	0.74
E	1.31	Q	1.07
$\mathbf{F}$	1.20	R	0.72
G	0.61	S	0.76
н	1.00	Т	0.89
I	1.32	V	1.25
K	0.96	Y	1.00
L	1.21	W	1.23

The penalty to be assigned to the gap depends on the type of residue, such as observed in sequences of known structure (Pascarella & Argos)

The penalty also depends on the location. If there are gaps around, the penalty increases

http://www.ebi.ac.uk/clustalw/ o to download the software

# **CLUSTALW** improvement

Thompson et al 1994





# CLUSTALW improvement

ClustalW is the most famous and the most used software for the progressive multi-alignment

It is available on EBI.

**ClustalX** is a new version of ClustalW that provide an intuitive GUI.

#### **Characteristics**

-The algorithm weights differently the sequences

and progressive changes matrix when add in the alignment proteins less similar.

-It is possible to add sequences to the ones already aligned, using pre-existing alignments.

-It is possible to perform structural alignments

Sequence WINDOW LENGTH det	SCORE TYPE		PAIRGAP
UNDOW LENGTH def	SCORE TYPE	TOPDIAG	PAIRGAP
def M OAP OPEN	percent ⊻	def 🛩	Local Inc.
OAP OPEN	110 00.000		cet 💌
	GAPS	OAP EXTENSION	OAP DISTANCES
def 💌	yes 🛩	det 💌	def 💌
ITERATION		NUMITE	R
none 👻		1 🛩	
1	PHYLO	GENETIC TREE	
OUTPUT TR ORDER	EETYPE CORRECT D	ST. IGNORE GAPS	5 CLUSTERING
aligned 😪 🛛 no	one 🖌 of 💌	off 💌	NJ
sequences in any s	upported format:		Help
			p
	def V ITERATION none V ORDER aligned In sequences in any s	ITERATION none PHYLO OUTPUT TREE TYPE CORRECT DI ORDER aligned Inone I of I requences in any supported format: Stoglia_	def     yes     def       ITERATION     NUMITE       none     1       OUTPUT     TREE TYPE       OUTPUT     TREE TYPE       ORDER     aligned       aligned     none       of     off       sequences in any supported format:

## CLUSTALW vs CLUSTALX

-------VHLTPEEKSAVTALWGKVN--VDEVGGEALGRLLVVYFWTQRFFESFGDLST ------VQLSGEEKAAVLALWDKVN--EEEVGGEALGRLLVVYFWTQRFFDSFGDLSN -----VLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFFTTKTYFPHFDLS-------VLSAADKTNVKAAWSKVGGHAGEYGAEALERMFLGFFTTKTYFPHFDLS-------VLSEGEWQLVLHVWAKVEADVAGHGQDILIRLFKSHFETLEKFDRFKHLKT PIVDTGSVAPLSAAEKTKIRSAWAPVYSTYETSGVDILVKFFTSTFAAQEFFPKFKGLTT -----GALTESQAALVKSSWEEFNANIPKHTHRFFILVLETAFAAKDLFSFLKGTSE

PDAVMGNPKVKAHGKKVLGAFSDGLAHLD----NLKGTFATLSELHCDKLHVDPENFRL PGAVMGNPKVKAHGKKVLHSFGEGVHHLD----NLKGTFAALSELHCDKLHVDPENFRL ---HGSAQVKGHGKKVADALTNAVAHVD----DMPNALSALSDLHAHKLRVDPVNFKL EAEMKASEDLKKHGVTVLTALGAILKKKG----HHEAELKPLAQSHATKHKIPIKYLEF ADQLKKSADVRWHAERIINAVNDAVASMDDT--EKMSMKLRDLSGKHAKSFQVDPQYFKV VP--QNNPELQAHAGKVFKLVYEAAIQLQVTGVVTDATLKNLGSVHVSKG-VADAHFPV CLUSTALW

LGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVANALAHKYH-----LGNVLVVVLARHFGKDFTPELQASYQKVVAGVANALAHKYH-----LSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTSKYR-----LSHCLLSTLAVHLPNDFTPAVHASLDKFLSSVSTVLTSKYR-----ISEAIIHVLHSRHPGDFGADAQGAMNKALELFRKDIAAKYKELGYQG LAAVIADTVAAG-----DAGFEKLMSMICILLRSAY-----VKEAILKTIKEVWGAKWSEELNSAWTIAYDELAIVIKKEMNDAA---

CLUSTALX.



#### Colors in the alignment

The raw-text files can be used to display the columns, but you can associate different colors for different physical residues with chemical characteristics. This greatly facilitates the display of multi-alignment.

7	G	W	QG <mark>V</mark>	WY I	R <mark>GWT</mark>	7ETA	K <mark>GL</mark> G	- <mark>LAG</mark>	WV <mark>R</mark> N	I <mark>R</mark> ADG	T – V	'EALF	<mark>h</mark> gp <b>e</b> ·	- <mark>AAV</mark> I	RAMLI	AC <mark>R</mark> G	- <mark>G</mark> -	-PPS <mark>A</mark>	R <mark>V</mark> DD <mark>L</mark>	RVT P	- <mark>VAAP</mark> -
1	G	UV	QG <mark>V</mark>	CY	R <mark>NW</mark> T (	7ENA	E <mark>QL</mark> G-	- <mark>IR</mark> G	WV <mark>R</mark> N	I <mark>rr</mark> dG	S− <mark>⊽</mark>	T <mark>EALF</mark>	SGPP	- E <mark>av</mark> i	) E <mark>mh</mark> q	R <mark>C</mark> RR	- <mark>G</mark> -	- PP <mark>AAI</mark>	MVTG <mark>l</mark>	E <mark>AF</mark> P	- <mark>st</mark> ee-
	G	V	QG <mark>V</mark>	GY	R <mark>AAC</mark> A	D <mark>MA</mark>	R <mark>AL</mark> G	- <mark>LR</mark> G	WV <mark>R</mark> N	I <mark>RR</mark> DG.	A-V	T <mark>EAFL</mark>	AGP E-	– PN <mark>VI</mark>	R <mark>NQ</mark> A	WMEE-	- <mark>G</mark> -	- PD LA	LVTQI	RTT P	GD <mark>IEP</mark> -
	GJ	CV	QG <mark>V</mark>	'GF I	RH <mark>AT</mark> I	7 <mark>RQA</mark>	I <mark>AL</mark> G	- <mark>IK</mark> G	WVA <mark>N</mark>	I <mark>L</mark> DDG	S-V	EAML	QGS <mark>A</mark> -	- <mark>NQV</mark> I	) <mark>R</mark> MLS	WL RH	- <mark>G</mark> -	- PP <mark>AA</mark>	R <mark>VTE</mark> V	'SG <mark>ee</mark>	RSTE <mark>R</mark> -
	GF	W	QG <mark>V</mark>	FF	R <mark>QS</mark> MF	EVA	H <mark>RNG</mark>	- <mark>VK</mark> G	WV <mark>R</mark> N	I <mark>R</mark> SDG	K T V	'EAVL	EGP <mark>R</mark> ·	-D <mark>AVI</mark>	<mark>1K</mark> VLE	WA <mark>R</mark> I <sup>.</sup>	- <mark>G</mark> -	- PPG <mark>A</mark>	R <mark>ved</mark> i	EVQW	EEY <mark>k</mark> G-
	G	NV.	QG <mark>V</mark>	SF	R <mark>ayt</mark> i	RD RA	R <mark>EAQ</mark> ·	- <mark>VR</mark> G	WV <mark>R</mark> N	ILSDG	R-V	'EAVF	EGT <mark>R</mark> -	- <mark>PAV</mark> (	2 <mark>K</mark> LIS	WCYS-	- <mark>G</mark> -	- <mark>PSQA</mark>	QVERV	'EV <mark>H</mark> W	EEPTG-
	G	NV.	QG <mark>V</mark>	GF	RWSM(	PEA	RK <mark>LG</mark> -	- <mark>VNG</mark>	WV <mark>R</mark> N	ILPDG	S-V	TEAVL	E <mark>g</mark> e e-	– E <mark>rv</mark> i	CALIG	WA <mark>H</mark> Q-	- <mark>C</mark> -	- PPFA	RVT RV	EV <mark>R</mark> W	EEP <mark>H</mark> G-
	G	W	QG <mark>V</mark>	GY	RYST (	7D T A	<mark>R</mark> QLG	- <mark>LTG</mark>	WV RN	IL PDN	R-V	'EAVF	EGA <mark>R</mark> ·	– E <mark>VV</mark> I	)D <mark>MV</mark> R	WC <mark>H</mark> S-	- <mark>G</mark> -	- PPAA	an <mark>k</mark> da	'VV <mark>E</mark> Y	E <mark>VPEG</mark> -
	G	V	QG <mark>V</mark>	GF	RYST (	ID T A	SQLC.	- <mark>LTG</mark>	WV RN	IL PD G	R-V	'EAVF	EG <mark>V</mark> R·	-D <mark>IV</mark> H	ED <mark>MV</mark> R	WC <mark>H</mark> A·	- <mark>G</mark> -	- PPAA	77 QD V	'AVEY	EE <mark>PEG</mark> -
2	Gl	ΓV	QGV	FF	RASMI	EEA.	L <mark>R</mark> LG	- <mark>LS</mark> G	WV <mark>R</mark> N	IL PD G	ESV	'EAVV	EG <mark>R</mark> G-	-D <mark>AV</mark> I	8 <mark>R</mark> IIC	WC L <mark>R</mark> ·	- <mark>G</mark> -	- PPAA	RVREL	RVEL	EPY <mark>k</mark> G-
	7	NV I	QGI	NF	RSNTI	' <mark>SK</mark> A	LELN.	-VRG	WVMN	I <mark>R</mark> IDG	s-v	'EAYF	SGE <mark>l</mark> ·	-CDVI	ISLIN	YCVS-	-D-	- <mark>MPYA</mark> Y	JVKR <mark>Y</mark>	DV-Y	D I PYM{
	G	W	QGI	NF	RSNTI	JV <mark>R</mark> A	LELG	-V <mark>R</mark> G	MIRN	ILPDG	S-V	TEALF	SGES	-EQI	<b>K</b> LIS	YCVS-	-N-	-MPYA	EVKRY	DV-Y	IEPYT-
	GI	W	QGV	CY	RQGT2	LQA	ERLA	- <mark>LA</mark> G	WV RN	ILADG	R-V	'EAWV	E <mark>c</mark> ee.	-AAVI	ELAE	MLM <mark>B</mark> .	-G-	- PEQA	RVEGV	ELEE	ACTOC-
	3	W	QG <mark>V</mark>	GF	QAT	EEA	D <mark>R</mark> LE	- <mark>L</mark> DC	WV <mark>R</mark> M	ILDDG	R-V	<mark>evvw</mark>	E <mark>c</mark> ee	-D <b>R</b> AF	<mark>lal e</mark> r	WLC <mark>R</mark>	- <mark>G</mark> -	- P <mark>RH</mark> A	evsav	'EVEQ	MPLQC;





#### **Using Colors**

### Colorazione

ALCDYNVFRDDIQPHWEVPENSNGGRWLIVIDKGKTPEMVDAIWLEILMALVGEQFGKDMESICGLVCNVRGKGSKISVW DLCDYNVFRDDIQPKWEAPENWDGGRWLIIINKGKTPEVLDAVWLEILLALIGEQFGKDMESICGLVCNVRGQGSKISVW DLCDYNVFRDDIQPKWEAPENWDGGRWLIIINKGKTPEVLDAVWLEILLALIGEQFGKDMESICGLVCNVRGQGSKISVW PPSDYNVFRDGIEPHWEVPQNQNGGRWLITIEKGRTPEIMDTIWTEILMAHIGEQFSDDIESLCGIVCNVRGKGSKISVW WGSDYYLFKEGIKPHWEDVNNVQGGRWLVVVDKQRRTQLLDHYWLELLMAIGEQFDEIGDYICGAVVNVRQKGDKVSLW SGCDYSLFKDGIEPHWEDERNRRGGRWLITLSKHQRKMDLDRFWLETLLCLVGEAFDDHSDDVCGAVINIRAKGDKIAIW SGCDYSLFKNGIQPKWEDAQNKKGGRWLINLNKTQRQTHLDDFWLETLLCLIGEGFDEHSEEICGATVNIRNKGDKLGLW LGSDYSLFKKNIRPHWEDAANKQGGRWVITLNKSSKTDL-DNLWLDVLLCLIGEAFDHS-DQICGAVINIRGKSNKISIW

Cysteine	C						
Negative	D,	Е					
Positive	Κ,	R					
Alcohol	S,	т					
Polar	N,	Q					
Aromatic	F,	Н,	W,	Y			
Hydrophobic	A,	G,	I,	L,	Μ,	Ρ,	V

## Raw Text files

				*	20	*	40	*	60			
	TRY1_HUMAN	:	IV <mark>GG</mark> YNCEE	NSVPYQV	SL <mark>N</mark>	SGYHFC	GGSLINEQWV:	/VS <mark>A</mark> GHCYKSR	IQVRLGEHNIEVLI	EGNE	:	62
1	TRYP_PIG	:	IVGGYTCAA	NSIPY <mark>Q</mark> V	SL <mark>N</mark>	SGSHFC	GGSLINSQWV:	/VS <mark>AAHC</mark> YKS <mark>R</mark>	IQVRLGEHNIDVLI	EGNE	:	62
	TRYU_DROME	:	IVGGADTSS	YYTKYV <mark>V</mark>	QLRRRS	SSSS <mark>S</mark> YAQT	GG <sup>CILDAVTI</sup>	at <mark>aahc</mark> vyn <mark>r</mark>	EAENFLVVSGDDSI	rggm	:	69
	TRYI_DROME	:	IIGGSDQLI	RNAPWQV	s <mark>IQ</mark>	ISARHEC	GG <mark>VIYSKEII:</mark>	ITAG <mark>HC</mark> LHER	SVT-LMKVRVGAQI	NHNY	:	62
	TRY1_SALSA	:	IV <mark>GG</mark> YECKA	YSQTHQV	SL <mark>N</mark>	SGYHFC	GGSLVNENWV	/VS <mark>AAHC</mark> YKS <mark>R</mark>	VEVRLGEHNIKVTI	EGSE	:	62
A												

			*	80	*	100	*	120	*	1		
7	TRY1_HUMAN	:	QFINAAK	IIR <mark>H</mark> PQY	DRKTLNNDI	MLIKLSSRAV	INARVSTI	SLPTAPPAT <mark>G</mark> T	kclis <mark>gwg</mark> n	TASS	:	127
2	TRYP_PIG	:	QFINAAK	IIT <mark>H</mark> PNFI	NGNTLDNDI	MLIKLSSPAT	LNSRVATV	SLPRSCAAA <mark>G</mark> T	eclis <mark>gwg</mark> n	TKSS	:	127
	TRYU_DROME	:	YGVVVRVSQ	LIP <mark>H</mark> ELY	NSSTMDNDI	ALVVVDPPLP	LDSFSTMEAT	VIASEQPPV <mark>G</mark> V	QATIS <mark>GWG</mark> Y	TKEN	:	138
6	TRYI_DROME	:	GGTLVPVAA	YKV <mark>h</mark> eqf	DSRFLHY <mark>DI</mark>	AVLRLSTPLT	FGLSTRAI	NLASTSPSG <mark>G</mark> T	TVTVT <mark>GWG</mark> H	TDNG	:	129
	TRY1_SALSA	:	QFISSSR	VIRHPNY	SSYNIDN <mark>DI</mark>	MLIKLSKPAT	LNTYVQPV	alptscapa <mark>g</mark> t	MCTVS <mark>GWG</mark> N	TMSS	:	127

			40	*	160	*	180	*	200		
	TRY1_HUMAN	:	GADYPDEL	CLDAP	/LSQAK <mark>C</mark> EASYPG	KITSN	M <mark>FC</mark> VG <mark>FL</mark> EGG	KDSCQGDSGGPV	/VCNGQLQ <mark>GVVS</mark>	:	192
6	TRYP_PIG	:	GSSYPSL <mark>L</mark>	CLKAP	/LSDSS <mark>C</mark> KSSYPG	QITGN	M <mark>IC</mark> VG <mark>FL</mark> EGG	KD <mark>SCQGDSGGP</mark> V	/VCNGQLQ <mark>G</mark> IVS	:	192
	TRYU_DROME	:	GLSS-DQ <mark>L</mark>	QVKVPJ	IVDSEK <mark>C</mark> QEAYYW	RPISEG	M <mark>lC</mark> aglsegg	KD <mark>ACQGDSGGP</mark> I	VVANKLAGIVS	:	203
	TRYI_DROME	:	ALSDS <mark>L</mark>	KAQLQI	LIDRGE <mark>C</mark> ASQKFG	YGADFVGEE	TI <mark>C</mark> AASTD	ADACTGDSGGPI	JVASSQLVGIVS	:	194
1	TRY1_SALSA	:	TADS-NKL	CLNI PJ	LSYSD <mark>C</mark> NNSYPG	MITNA	M <mark>FC</mark> AGYLEGG	KD <mark>SCQGDSGGP</mark> V	VCNGELQGVVS	:	191
2											

			*	220	*		
	TRY1_HUMAN	:	WGDGCAQ	KNK <mark>PGVY</mark> TKV	YNYVK <mark>WI</mark> KNTI.	AANS :	224
	TRYP_PIG	:	WGYGCAQ	KNKPGVY <mark>TK</mark> V	CNYVNWIQQTI.	AAN- :	223
	TRYU_DROME	:	WGEGCAR	PNY <mark>PGVY</mark> ANV	AYYKDWIAKQR	TSYV :	235
	TRYI_DROME	:	WGYR <mark>CA</mark> D	DNY PGVY ADV	AILRPWIVKAA	NAI- :	225
6	TRY1_SALSA	:	WGYGCAE	pgn <mark>pgvy</mark> ak <mark>v</mark>	CIFNDWLTSTM	ASY- :	222

## Raw text files

		*	20	*	40	*	60		
TRY1_HUMAN	:	IVGGYNCEENSVPY	QVSLN	-SGYHFCG	GSLINEQWVVSA	GHCYKSRIÇ	VRLGEHNIEVLEGNE	:	62
TRYP_PIG	:	T.AAI		s	s	А	D	:	62
TRYU_DROME	:	ADTSSYYTK.	V.Q.RRRSSS	SS.YAQT	.CILDAVTIAT	AVYN.EA	ENFLVVSGDDSR.GM	:	69
TRYI_DROME	:	.ISDQLIRNA.W	IQ	I.AR.E	.VIYSKEIIIT.	LHE.SV	TMKVRVGAQNH.Y	:	62
TRY1_SALSA	:	E.KAY.QTH			VN	AVE	K.TS.	:	62

		* 80	*	100	*	120	* 1		
TRY1_HUMAN	:	QFINAAKIIRHPQ	YDRKTLNNDI	MLIKLSSRAVI	ENARVSTI	SLPTAPPATG	TKCLISGWGNTASS	:	127
TRYP_PIG	:		FNGND	P.T]	LSA.V	RSCA.A.	.EK	:	127
TRYU_DROME	:	YGVVVRVSQL.P.EL	.NSS.MD	A.VVVDPPLP	LDSFSTMEA.	VIASEQ.PV.	VQATY.KEN	:	138
TRYI_DROME	:	GGTLVPV.AYKV.E.	F.SRF.HY	AVLRTPLTH	GLSTRA.	N.ASTS.SG.	.TVTVTH.DNG	:	129
TRY1 SALSA	:	N	.SSYNID	KP.TI	LTY.QPV.	ASCAPA.	.M.TVM	:	127

		40 *	160	*	180	*	200		
TRY1_HUMAN	:	GADYPDELQCLDAE	VLSQAKCEASYPG-	KITSNI	MFCVGFLEGG	KDSCQGDSGG	PVVCNGQLQGVVS	:	192
TRYP_PIG	:	.SSSLK	DSS.KS	QG.	. I		I	:	192
TRYU_DROME	:	.LSSQQVKV.	IVDSEQEA.YWF	RP.SEG	.L.A.LS	A	.L.VANK.A.I	:	203
TRYI_DROME	:	ALS.SKAQLQ	DIDRGE.ASQKF.N	'GADEVGEE'	TI.AASTDA	Α.Α.Τ	.L.ASSV.I	:	194
TRY1_SALSA	:	TS-NKNI.	IYSD.NN	MNA	A.Y		E	:	191

		* 220 *		
TRY1_HUMAN	:	WGDGCAQKNKPGVYTKVYNYVKWIKNTIAANS	:	224
TRYP_PIG	:	YCNQQ	:	223
TRYU_DROME	:	ERP.YAN.AY.KDAKQRTSYV	:	235
TRYI_DROME	:	YRDD.YAD.AILRPVKAAN.I-	:	225
TRY1_SALSA	:	YEPGNACIFND.LTS.M.SY-	:	222





#### **Consensus sequences**

A consensus sequence is defined as a sequence derived by a multi-alignment that presets only most conserved residues for each position:

- $\Rightarrow$  it **summarizes** a multi-alignment.
- $\Rightarrow$  it **is different** to all proteins of the input dataset.

 $\Rightarrow$  it is possible to define special symbols that specify not perfect conservation in some position.

 $\Rightarrow$  it is possible to use a specific syntax that also allows to specify variants in a position, other than conservations.



### **Consensus Sequences**

To facilitate the reading of a multiple alignment and quickly highlight conserved amino acids, it can be represented the "consesus" sequence in the last line of a multi-alignment.

If all the sequences in a multiple alignment contains the same amino acid in a location than it is represented in the consensus sequence, in the same way if all the amino acids of a column belong to the same class, in the consensus line it is shown the symbol of class.

It can also be represented below the consensus sequence, the rows of the consensus that below a certain percentage (for example 90%, 80% or 70%). In these lines, to put a symbol it is sufficient that the amino acid is conserved in a population of amino acids of the upper column to the percentage of the consensus (rather than in all).







PEEKSAVTALWGKVN--VDEVGG Hbb\_human GEEKAAVLALWDKVN--EEEVGG Hbb\_horse PADKTNVKAAWGKVGAHAGEYGA Hba\_human AADKTNVKAAWSKVGGHAGEYGA Hba\_horse EHEWQLVLHVWAKVEAGVAGHGQ Myg whale