

Chapter III: Pairwise Sequence Alignment

Presentations use info from:

Jonathan Pevsner, Ph.D.
<http://bioinfbook.org>
pevsner@kennedykrieger.org
Bioinformatics and Functional Genomics
(3rd edition, ©2015 John Wiley & Sons, Ltd.)
You may use this PowerPoint for teaching purposes

- dr. Stanislav Kolenčík
stanislav.kolencik@famnit.upr.si

What will you learn?

To define homology as well as orthologs and paralogs

To explain how PAM (accepted point mutation) matrices are derived

To contrast the utility of PAM and BLOSUM scoring matrices.

To define dynamic programming and explain how global and local pairwise alignments are performed

To perform pairwise alignment of protein or DNA sequences at the NCBI website

Pairwise sequence alignment is the most fundamental operation of bioinformatics

- It is used to decide if two proteins (or genes) are related structurally or functionally.
- It is used to identify domains or motifs that are shared between proteins.
- It is the basis of BLAST searching.
- It is used in the analysis of genomes.

Pairwise alignment is the process of lining up two sequences to achieve maximal levels of identity (and maximal levels of conservation in the case of amino acid alignments) for the purpose of assessing the degree of similarity and the possibility of homology.

A glowing blue DNA double helix structure is visible on the right side of the image, set against a dark blue background with a subtle light gradient.

Protein alignment: often
more informative than DNA
alignment

Sequence alignment: protein sequences can be more informative than DNA

Protein is more informative (20 vs 4 characters);

Many amino acids share related biophysical properties

Codons are degenerate: changes in the third position

Often do not alter the amino acid that is specified

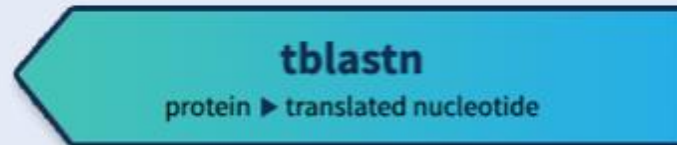
Protein sequences offer a longer “look-back” time

Sequence alignment: protein sequences can be more informative than DNA

Example:

- >searching for plant globins using human beta globin DNA yields no matches;
- >searching for plant globins using human beta globin protein yields many matches.

Web BLAST



Pairwise sequence alignment is the most fundamental operation of bioinformatics

Many times, DNA alignments are appropriate:

- to study non-coding regions of DNA (e.g., introns or intergenic regions)
- database searching
- to study DNA polymorphisms
- genome sequencing relies on DNA analysis

Definitions

Homology

- Similarity attributed to descent from a common ancestor.

Identity

- The extent to which two (nucleotide or amino acid) sequences are invariant.

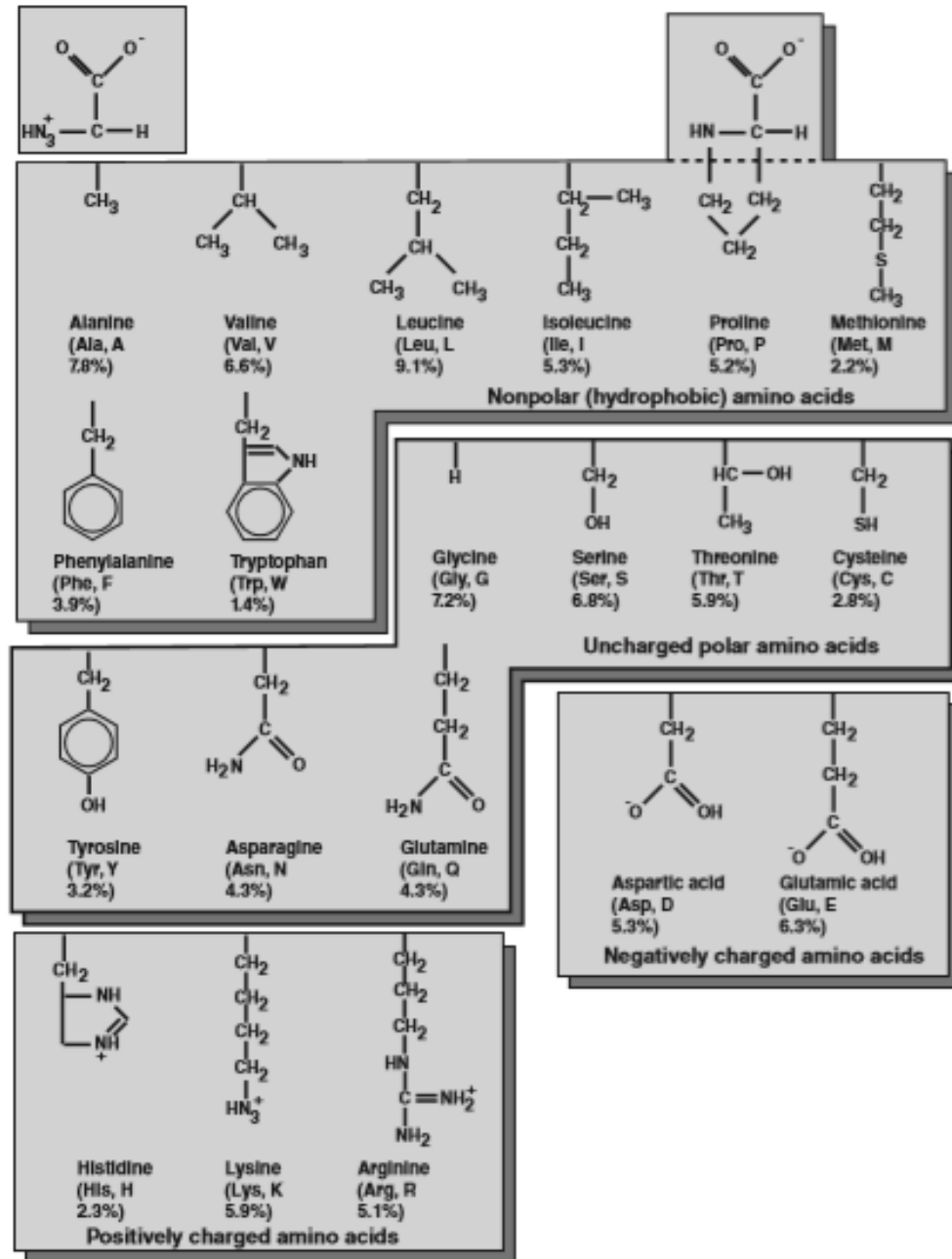
Similarity

- The extent to which nucleotide or protein sequences are related. It is based upon identity plus conservation.

Conservation

- Changes at a specific position of an amino acid or (less commonly, DNA) sequence that preserve the physio-chemical properties of the original residue.

BOX 3.2 STRUCTURES AND ONE- AND THREE-LETTER ABBREVIATIONS OF 20 COMMON AMINO ACIDS



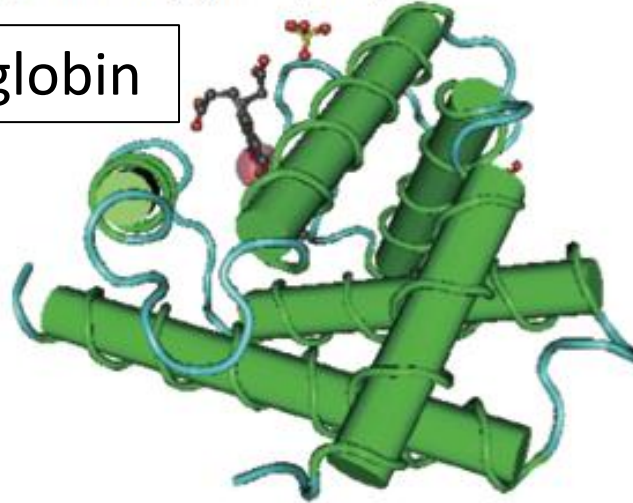
- basic amino acids (K, R, H)
- acidic amino acids (D, E)
- hydroxylated amino acids (S, T)
- hydrophobic amino acids (W, F, Y, L, I, V, M, A)

Globin homologs

In general, three-dimensional structures diverge much more slowly than amino acid sequence identity between two proteins.

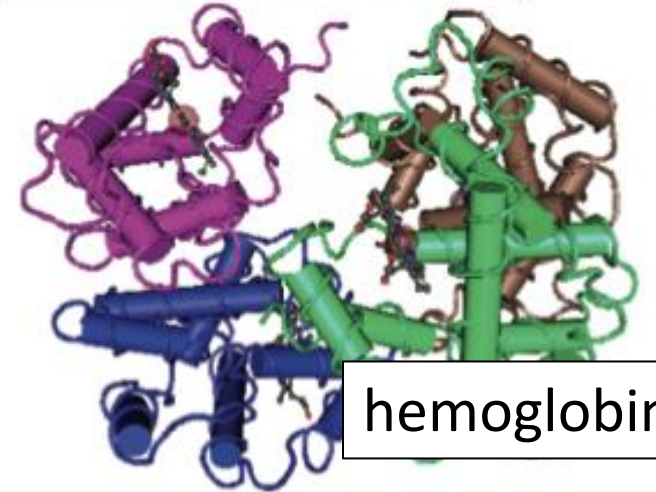
(a) Human myoglobin (3RGK)

myoglobin



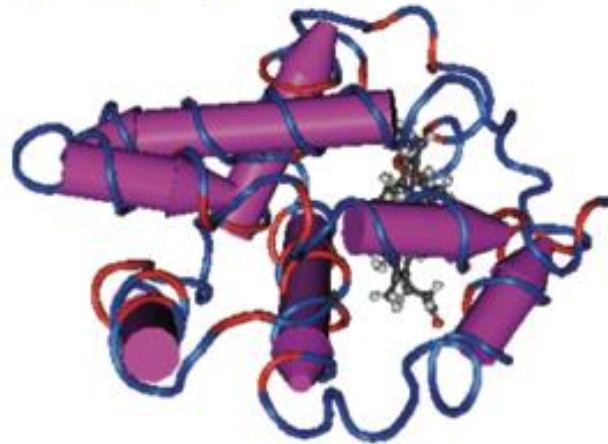
(b) Human hemoglobin tetramer (2H35)

hemoglobin

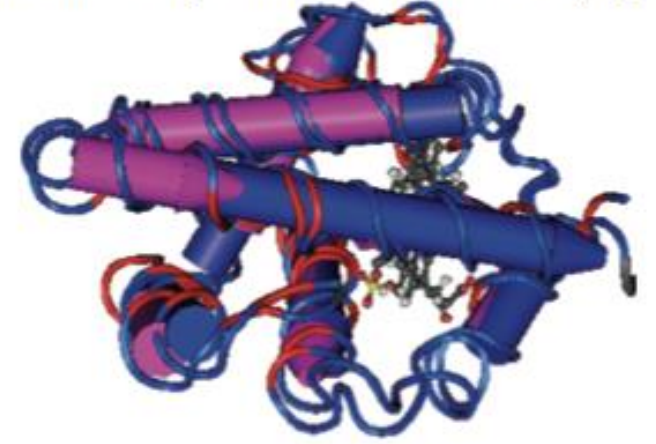


(c) Human beta globin (subunit of 2H35)

beta globin



(d) Pairwise alignment of beta globin and myoglobin



beta globin and myoglobin
(aligned)



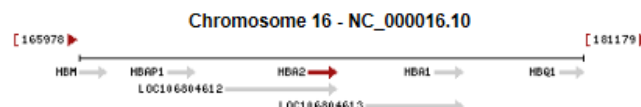
Definitions: two types of homology

Orthologs

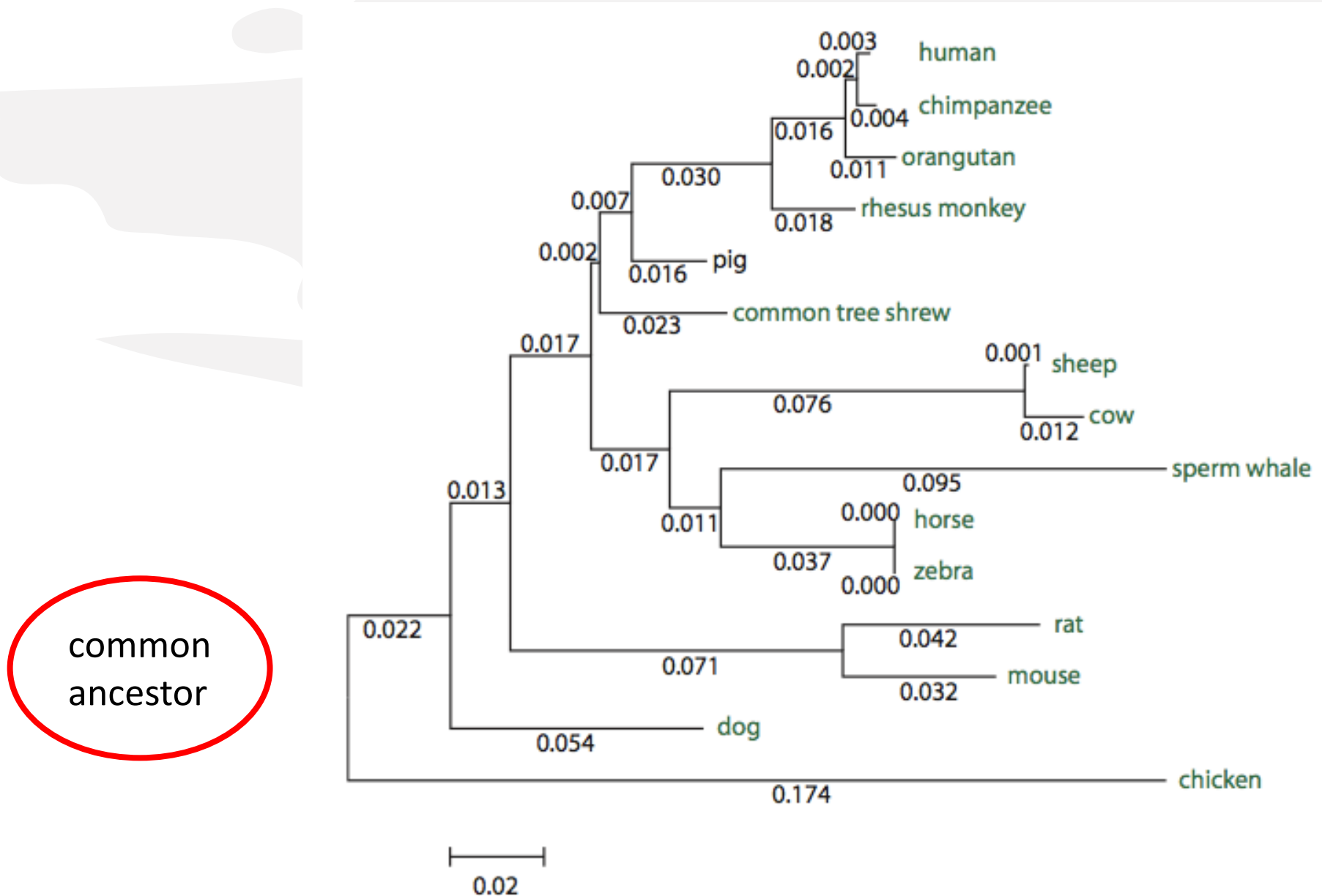
- Homologous sequences in different species that arose from a common ancestral gene during speciation; may or may not be responsible for a similar function.

Paralogs

- Homologous sequences within a single species that arose by gene duplication.

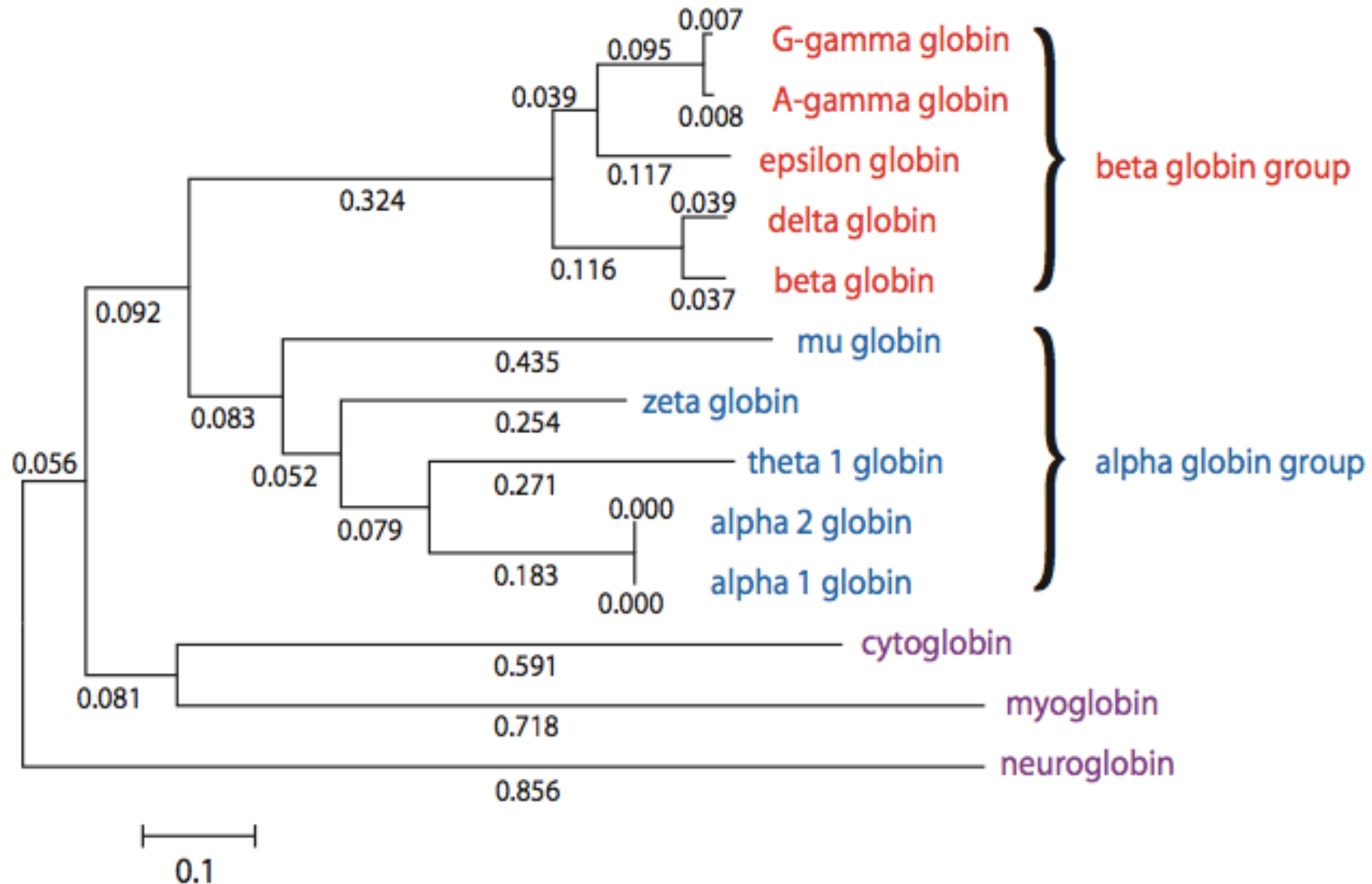


Myoglobin proteins: examples of orthologs



Paralogs: members of a gene (protein) family within a species

This tree shows human globin paralogs.

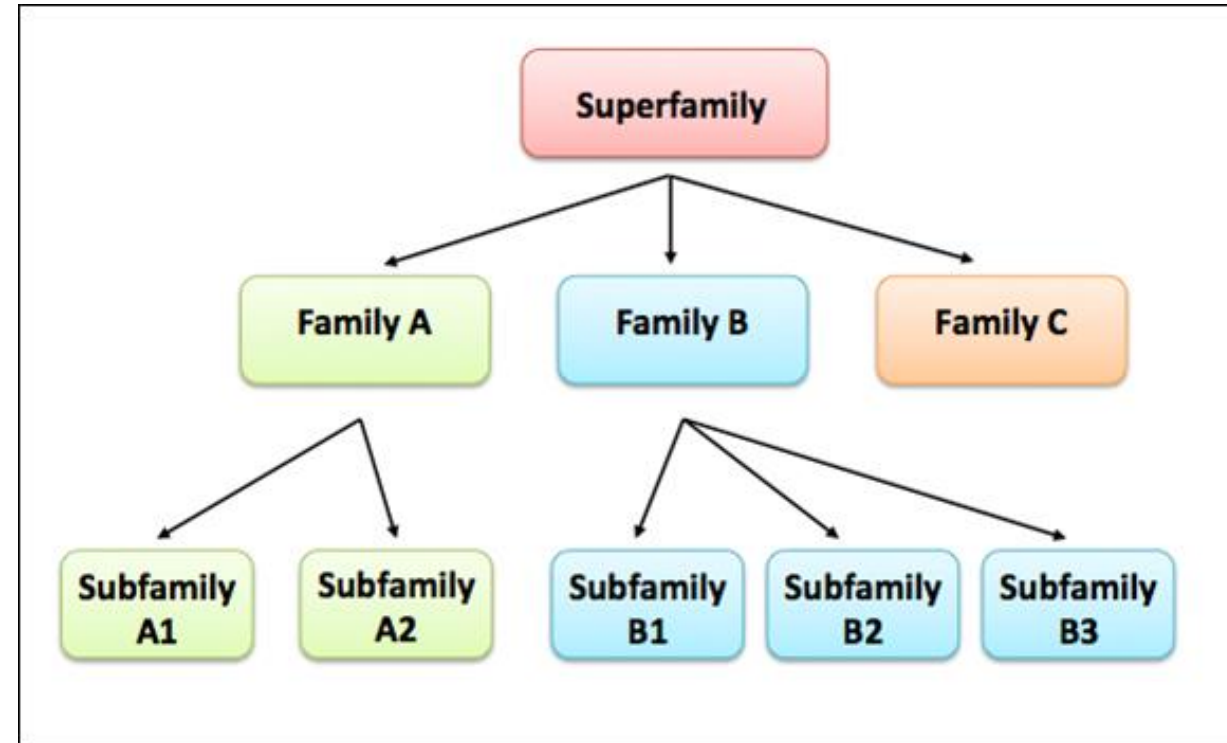


Why classify proteins?

- Proteins can be classified into groups according to sequence or structural similarity.
- These groups often contain well characterized proteins whose function is known.
- When a novel protein is identified, its functional properties can be proposed based on the group to which it is predicted to belong.

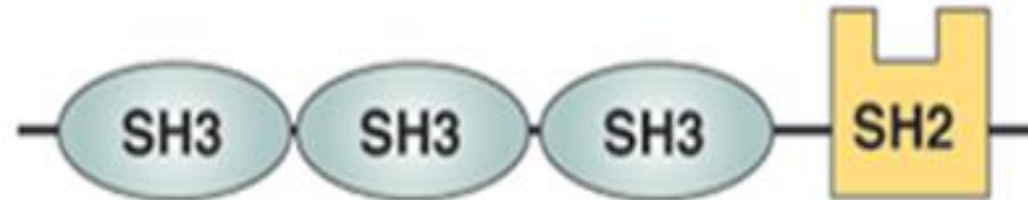
What are protein families?

A group of proteins that share a common evolutionary origin, reflected by their related functions and similarities in sequence or structure.



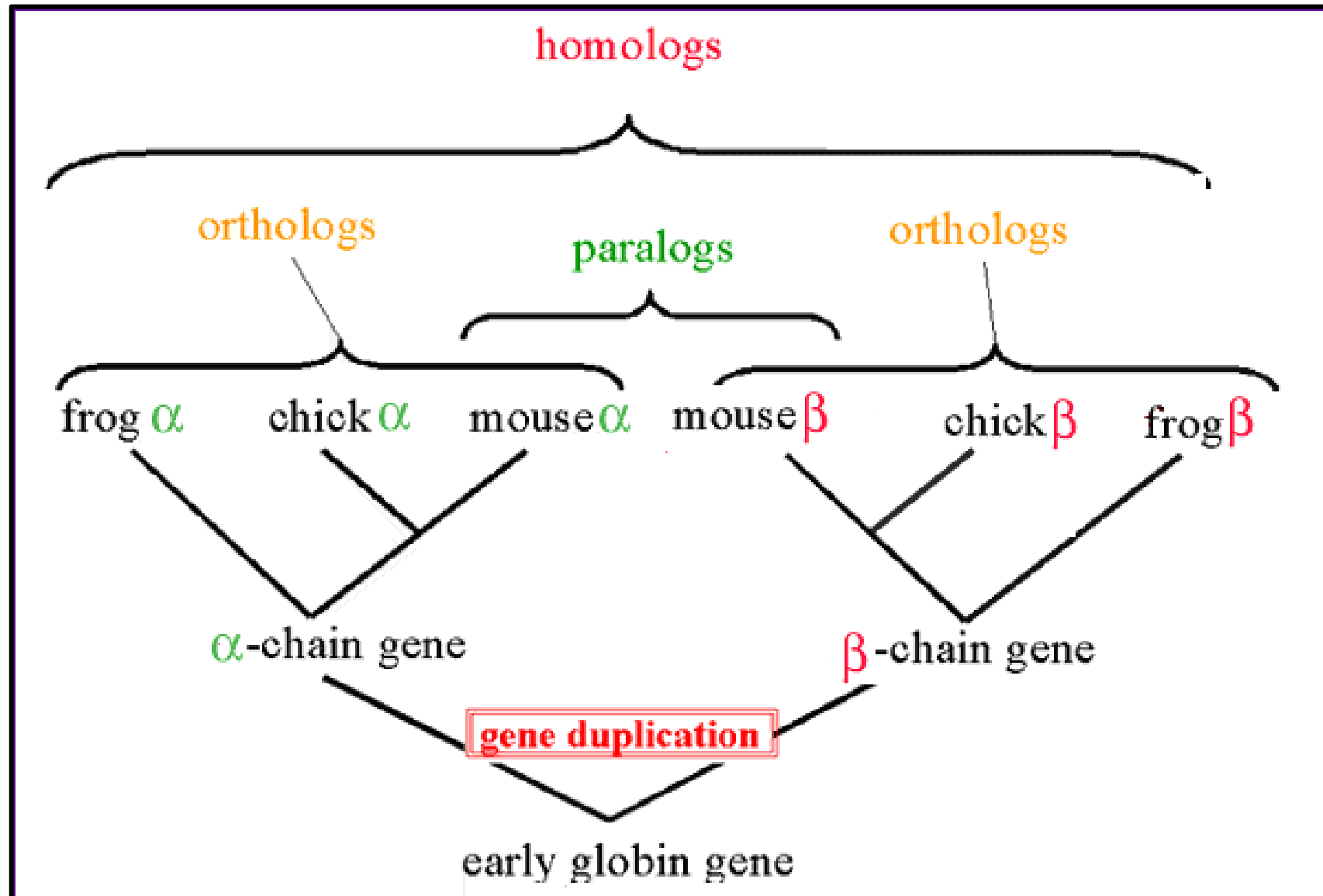
What are protein domains?

- Domains are distinct functional and/or structural units in a protein.
- They are responsible for a particular function or interaction, contributing to the overall role of a protein.
- Domains may exist in a variety of biological contexts, where similar domains can be found in proteins with different functions.



Domain composition of Nck

Orthologs and paralogs are often viewed in a single tree





General approach to pairwise alignment

- ✓ Choose two sequences
- ✓ Select an algorithm that generates a score
- ✓ Allow gaps (insertions, deletions)
- ✓ Score reflects degree of similarity
- ✓ Alignments can be global or local
- ✓ Estimate probability that the alignment occurred by chance

Find BLAST from the home page of NCBI and select protein BLAST...

Popular Resources

[PubMed](#)

[Bookshelf](#)

[PubMed Central](#)

[BLAST](#)

[Nucleotide](#)

[Genome](#)

[SNP](#)

[Gene](#)

[Protein](#)

[PubChem](#)

Basic Local Alignment Search Tool

BLAST finds regions of similarity between biological sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance. [Learn more](#)

The Basic Local Alignment Search Tool (BLAST) finds regions of local similarity between sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance of matches. BLAST can be used to infer functional and evolutionary relationships between sequences as well as help identify members of gene families.

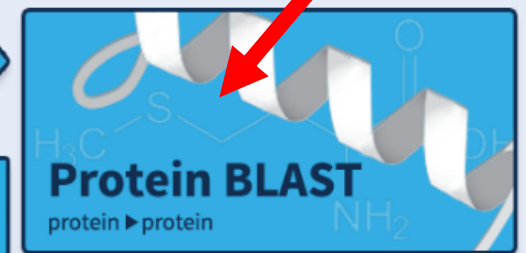
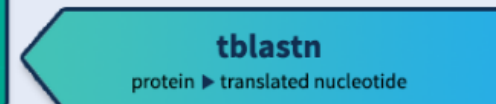
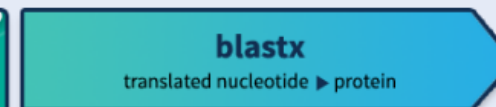
ClusteredNR database on BLAST+

The ClusteredNR database is now available for BLAST+

Thu, 24 Aug 2023

[More BLAST news...](#)

Web BLAST



- Choose align two or more sequences...

Bookmark

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) [?](#) [Clear](#)

Query subrange [?](#)

From

To

Or, upload file No file chosen [?](#)

Job Title

Enter a descriptive title for your BLAST search [?](#)

☐ Align two or more sequences [?](#)

Choose Search Set

Databases ☒ Standard databases (nr etc.): [New](#) ☐ Experimental databases

Compare ☐ Select to compare standard and experimental database [?](#)

Standard

Database [?](#)

Organism Optional ☐ exclude [Add organism](#)

Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown [?](#)

Exclude Optional ☐ Models (XM/XP) ☐ Non-redundant RefSeq proteins (WP) ☐ Uncultured/environmental sample sequences

Program Selection

Algorithm

☐ Quick BLASTP (Accelerated protein-protein BLAST)

☒ blastp (protein-protein BLAST)

☐ PSI-BLAST (Position-Specific Iterated BLAST)

☐ PHI-BLAST (Pattern Hit Initiated BLAST)

☐ DELTA-BLAST (Domain Enhanced Lookup Time Accelerated BLAST)

Choose a BLAST algorithm [?](#)

[BLAST](#) Search **database nr** using **Blastp (protein-protein BLAST)**

☐ Show results in a new window

[+ Algorithm parameters](#)

[Feedback](#)

blastn **blastp** blastx tblastn tblastx

BLASTP programs search protein subjects using a protein query. more...

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) ? Clear

NP_000509

Query subrange ?

From

To

Or, upload file

Choose File No file chosen ?

Job Title

NP_000509:hemoglobin subunit beta [Homo sapiens]

Enter a descriptive title for your BLAST search ?

☒ Align two or more sequences ?

Enter Subject Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) ? Clear

NP_005359

Subject subrange ?

From

To

Or, upload file

Choose File No file chosen ?

Program Selection

Algorithm

☒ blastp (protein-protein BLAST)

Choose a BLAST algorithm ?

BLAST

Search protein sequence using Blastp (protein-protein BLAST)

☐ Show results in a new window

+ Algorithm parameters

Enter the two sequences
(as accession numbers
or in the fasta format)
and click **BLAST**.

NP_000509
NP_005359

Optionally select “Algorithm parameters” and note the matrix option.

Algorithm parameters

General Parameters

Max target sequences 100 ?
Select the maximum number of aligned sequences to display ?

Short queries ☒ Automatically adjust parameters for short input sequences ?

Expect threshold 0.05 ?

Word size 3 ?

Max matches in a query range 0 ?

Scoring Parameters

Matrix BLOSUM45 ?

Gap Costs Existence: 13 Extension: 3 ?

Compositional adjustments Conditional compositional score matrix adjustment ?

Filters and Masking

Filter ☐ Low complexity regions ?

Mask ☐ Mask for lookup table only ?
☐ Mask lower case letters ?

BLAST

Search protein sequence using Blastp (protein-protein BLAST)

☐ Show results in a new window

[blastn](#)**blastp**[blastx](#)[tblastn](#)[tblastx](#)BLASTP programs search protein subjects using a protein query. [more...](#)**Enter Query Sequence**Enter accession number(s), gi(s), or FASTA sequence(s) [?](#) [Clear](#)

>NP_000509.1 hemoglobin subunit beta [Homo sapiens]
MVHLTPEEKSAVTALWGKLVNDEVGGEALGRLLVYPWTQRFESFGDLST
PDAVMGNPKVKAHGKKVLG
AFSDGLAHLNLIKGTATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKE

Query subrange [?](#)

From

To

Or, upload file

Choose File

No file chosen

[?](#)

Job Title

Enter a descriptive title for your BLAST search [?](#)Align two or more sequences [?](#)**Enter Subject Sequence**Enter accession number(s), gi(s), or FASTA sequence(s) [?](#)[Clear](#)Subject subrange [?](#)

NP_005359

From

To

Or, upload file

Choose File

No file chosen

[?](#)**Program Selection**

Algorithm

☒ blastp (protein-protein BLAST)Choose a BLAST algorithm [?](#)**BLAST**Search **protein sequence** using **Blastp (protein-protein BLAST)**

Show results in a new window

+ Algorithm parameters

Pairwise alignment of human beta globin (the “query”) and myoglobin (the “subject”)

myoglobin isoform 1 [Homo sapiens]

Sequence ID: [NP_001349775.1](#) Length: 154 Number of Matches: 1

[See 12 more title\(s\)](#) ▼ [See all Identical Proteins\(IPG\)](#)

Range 1: 3 to 147 [GenPept](#) [Graphics](#)

▼ [Next Match](#) ▲ [Previous Match](#)

Score = 43.9 bits (102), Expect = 1e-09, Method: Composition-based stats.

Identities = 37/145 (25%), Positives = 57/145 (39%), Gaps = 2/145 (1%)

```

Query  4      LTPEEKSAVTALWGKVNVD--EVGGEALGRLLVVYPWTQRFFESFGDLSTPDVAVMGNPKV  61
      → L+  E   V  +WGKV  D    G E L RL   +P T   F+ F   L + D +   +   +
Sbjct  3      LSDGEWQLVLNVWGKVEADIPGHGQEVLIIRLFKGHHPETLEKFDKFKHLKSEDEMKASEDL  62

Query  62     KAHGKKVLGAFSDGLAHLNLTGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGK  121
      → K HG  VL A    L    + +      L++ H K  +   +   +   ++ VL
Sbjct  63     KKHGATVLTALGGILKKKGHHEAEIKPLAQSHATKHKIPVKYLEFISECIIQVLQSKHPG  122

Query  122    EFTPPVQAAYQKVVAGVANALAHKY  146
      → +F    Q A  K +      +A  Y
Sbjct  123    DFGADAQGAMNKALELFRKDMASNY  147
```

We'll examine the highlighted green region of the alignment in more detail.

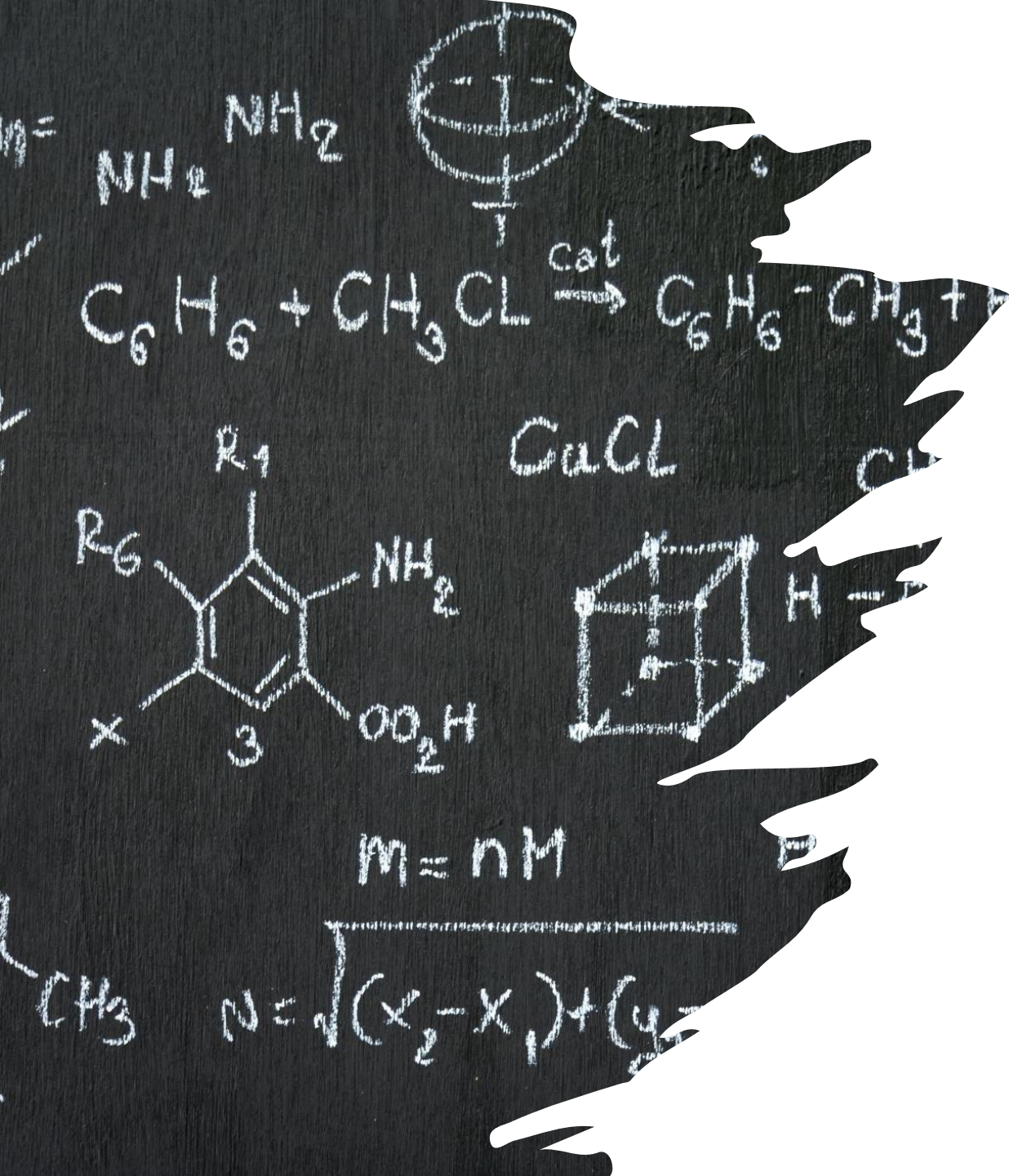
How raw scores are calculated: an example

```
Score = 18.1 bits (35), Expect = 0.015, Method: Composition-based stats.  
Identities = 11/24 (45%), Positives = 12/24 (50%), Gaps = 2/24 (8%)
```

```
Query 12 VTALWGKVNVD--EVGGEALGRLL 33  
          V +WGKV D G E L RL  
Sbjct 11 VLNVWGKVEADIPGHGQEVLRLF 34
```

match	4	11	5	6	6	5	4	5	sum of matches: +60 (round up to +61)
		6	4					4	
mismatch	-1	1	0	-2	-2	-4	0		sum of mismatches: -13
	-2		0	-3	0				
gap open			-11						sum of gap penalties: -13
gap extend			-2						
									total raw score: 61 - 13 - 13 = 35

For a set of aligned residues, we assign scores based on matches, mismatches, gap open penalties, and gap extension penalties. These scores add up to the total raw score.



Where do scores come from?

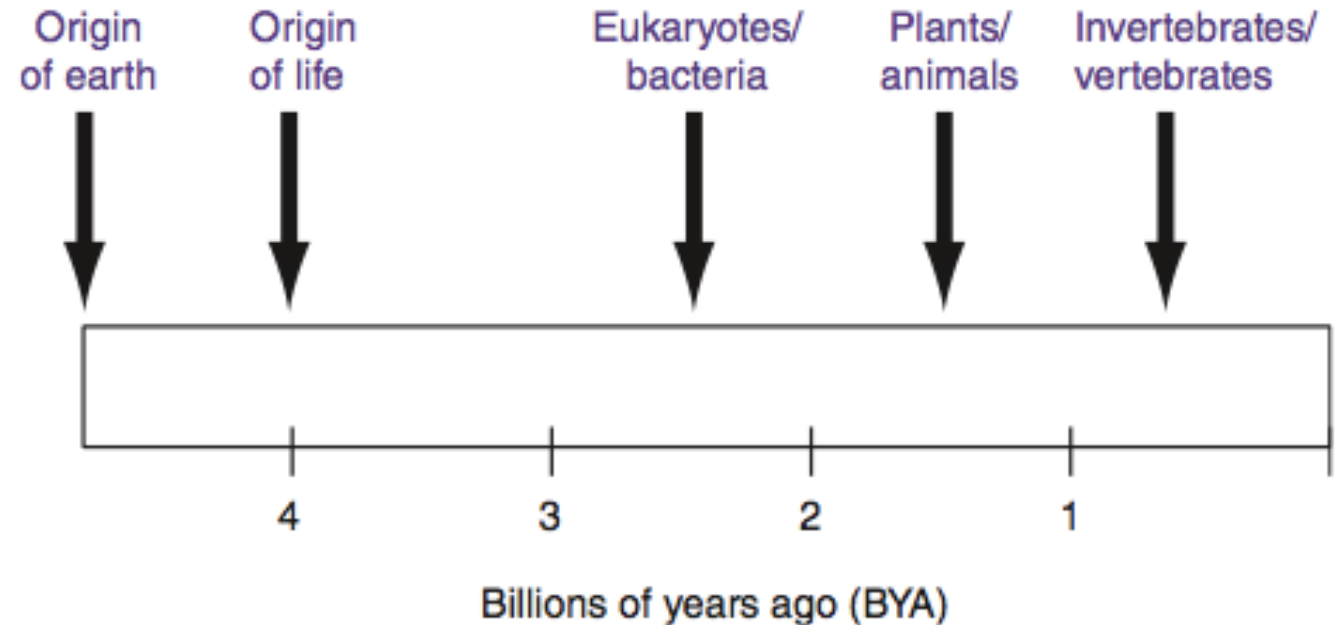
We'll examine scoring matrices. These are related to the properties of the 20 common amino acids.

Gap

- Positions at which a letter is paired with a null are called gaps.
- Gap scores are typically negative.
- Since a single mutational event may cause the insertion or deletion of more than one residue, the presence of a gap is ascribed more significance than the length of the gap. Thus, there are separate penalties for gap creation and gap extension.
- In BLAST, it is rarely necessary to change gap values from the default.

Pairwise alignment and the evolution of life

- When two proteins (or DNA sequences) are homologous they share a common ancestor.
- We can infer the sequence of that ancestor.
- When we align globins from human and a plant we can imagine their common ancestor, a single celled organism that lived 1.5 billion years ago, and we can infer that ancient globin sequence.
- Through pairwise alignment we can look back in time at sequence evolution.

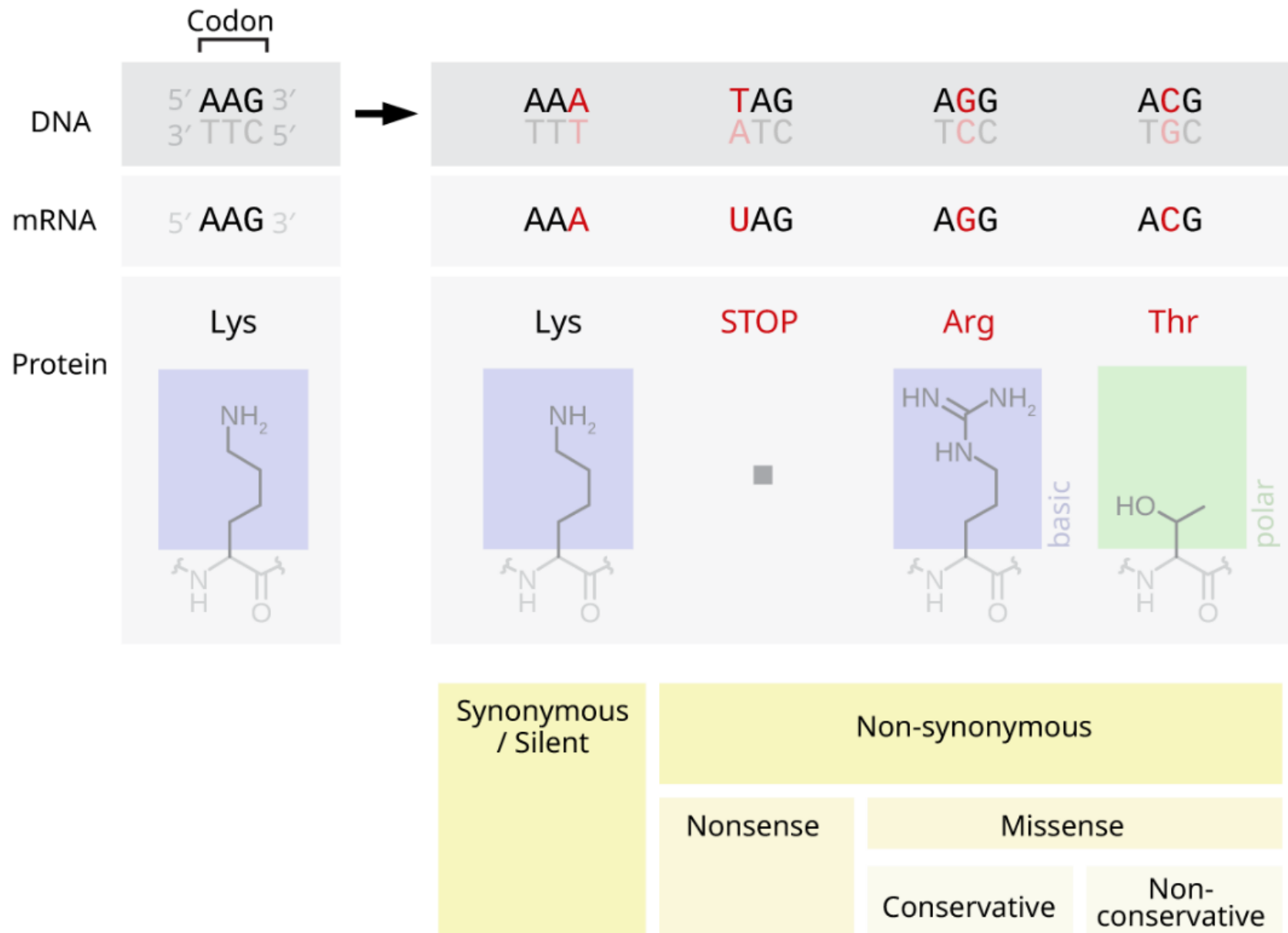


Scoring matrices -> Dayhoff model: 7 steps

Step 1: Accepted point mutations (PAMs) in protein families

PROTEIN	PAMS PER 100 MILLION YEARS
Immunoglobulin (Ig) kappa chain C region	37
Kappa casein	33
Epidermal growth factor	26
Serum albumin	19
Hemoglobin alpha chain	12
Myoglobin	8.9
Nerve growth factor	8.5
Trypsin	5.9
Insulin	4.4
Cytochrome c	2.2
Glutamate dehydrogenase	0.9
Histone H3	0.14
Histone H4	0.10

Margaret Dayhoff and colleagues developed scoring matrices in the 1960s and 1970s. They defined PAMs as “accepted point mutations.” Some protein families evolve very slowly (e.g. histones change little over 100 million years); others (such as kappa casein) change very rapidly.



An example of point mutations at an amino acid site coding for lysine. The missense mutations may be classed as **point accepted mutations** if the **mutated protein is not rejected by natural selection** (Source: Wikipedia).

Dayhoff's 34 protein superfamilies

<u>Protein</u>	<u>PAMs per 100 million years</u>
Ig kappa chain	37
Kappa casein	33
lutinizing hormone b	30
lactalbumin	27
complement component 3	27
epidermal growth factor	26
proopiomelanocortin	21
pancreatic ribonuclease	21
haptoglobin alpha	20
serum albumin	19
phospholipase A2, group IB	19
prolactin	17
carbonic anhydrase C	16
Hemoglobin a	12
Hemoglobin b	12

Dayhoff's 34 protein superfamilies

<u>Protein</u>	<u>PAMs per 100 million years</u>
Ig kappa chain	37
Kappa casein	33
luteinizing hormone b	30
lactalbumin	27
complement component 3	27
epidermal growth factor	26
proopiomelanocortin	21
pancreatic ribonuclease	21
haptoglobin alpha	20
serum albumin	19
phospholipase A2, group IB	19
prolactin	
carbonic anhydrase C	16
Hemoglobin a	12
Hemoglobin b	12

human (NP_005203) versus mouse (NP_031812) kappa casein

```
Score = 57.8 bits (138), Expect = 3e-07
Identities = 39/118 (33%), Positives = 61/118 (51%), Gaps = 2/118 (1%)

Query 1 MKSFLLVVNALALTLPFLAVEVQNQKQPACHENDERPFYQKTAPYVPMYYVPNSYPYYGT 60
M++F++V+N LALTLPFLA E+QN E ++ + ++ Y P+ V N + Y
Sbjct 2 MRNFIVVMNILALTLPFLAAEIQNPDSNCRGEKNDIVYDEQRVLYTPVRSVLN-FNQYEP 60

Query 61 NLYQRRPAI-AINNPFYVPRTYYANPAVVRPHAQIPQRQYLPNSHPPTVVRPLNLHPSF 117
N Y RP++ A +PY+ ++R A I + Q +PN V +PSF
Sbjct 61 NYYHYRPSLPATASPYMYYPVVRLLLLLRSPAPISKWQSMFPNPQSAGVPYAIPNPSF 118
```

Dayhoff's 34 protein superfamilies

<u>Protein</u>	<u>PAMs per 100 million years</u>
apolipoprotein A-II	10
lysozyme	9.8
gastrin	9.8
myoglobin	8.9
nerve growth factor	8.5
myelin basic protein	7.4
thyroid stimulating hormone b	7.4
parathyroid hormone	7.3
parvalbumin	7.0
trypsin	5.9
insulin	4.4
calcitonin	4.3
arginine vasopressin	3.6
adenylate kinase I	3.2

Dayhoff's 34 protein superfamilies

<u>Protein</u>	<u>PAMs per 100 million years</u>
triosephosphate isomerase I	2.8
vasoactive intestinal peptide	2.6
glyceraldehyde phosph. dehydrogease	2.2
cytochrome c	2.2
collagen	1.7
troponin C, skeletal muscle	1.5
alpha crystallin B chain	1.5
glucagon	1.2
glutamate dehydrogenase	0.9
histone H2B, member Q	0.9
ubiquitin	0

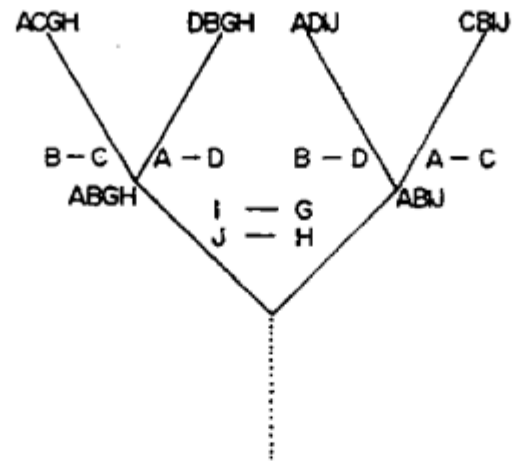


Figure 78. Simplified phylogenetic tree. Four "observed" proteins are shown at the top. Inferred ancestors are shown at the nodes. Amino acid exchanges are indicated along the branches.

	A	B	C	D	G	H	I	J
A			1	1				
B			1	1				
C	1	1						
D	1	1						
G							1	
H								1
I					1			
J						1		

Figure 79. Matrix of accepted point mutations derived from the tree of Figure 78.

Step 1: accepted point mutations are defined not by the pairwise alignment but with respect to the common ancestor



Dayhoff et al. evaluated amino acid changes. They applied an evolutionary model to compare changes such as 1 versus 2 not to each other but to an inferred common ancestor at position 5.

	A Ala	R Arg	N Asn	D Asp	C Cys	Q Gln	E Glu	G Gly	H His	I Ile	L Leu	K Lys	M Met	F Phe	P Pro	S Ser	T Thr	W Trp	Y Tyr	V Val
A																				
R	30																			
N	109	17																		
D	154	0	532																	
C	33	10	0	0																
Q	93	120	50	76	0															
E	266	0	94	831	0	422														
G	579	10	156	162	10	30	112													
H	21	103	226	43	10	243	23	10												
I	66	30	36	13	17	8	35	0	3											
L	95	17	37	0	y	75	15	17	40	253										
K	57	477	322	85	0	147	104	60	23	43	39									
M	29	17	0	0	0	20	7	7	0	57	207	90								
F	20	7	7	0	0	0	0	17	20	90	167	0	17							
P	345	67	27	10	10	93	40	49	50	7	43	43	4	7						
S	772	137	432	98	117	47	86	450	26	20	32	168	20	40	269					
T	590	20	169	57	10	37	31	50	14	129	52	200	28	10	73	696				
W	0	27	3	0	0	0	0	0	3	0	13	0	0	10	0	17	0			
Y	20	3	36	0	30	0	10	0	40	13	23	10	0	260	0	22	23	6		
V	365	20	13	17	33	27	37	97	30	661	303	17	77	10	50	43	186	0	17	
	A Ala	R Arg	N Asn	D Asp	C Cys	Q Gln	E Glu	G Gly	H His	I Ile	L Leu	K Lys	M Met	F Phe	P Pro	S Ser	T Thr	W Trp	Y Tyr	V Val

From a survey of 1572 observed substitutions, the original amino acid (columns) are compared to the changes (rows).

Dayhoff model step 2 (of 7): Frequency of amino acids

TABLE 3.1 Normalized frequencies of amino acid. These values sum to 1. If the 20 amino acids were equally represented in proteins, these values would all be 0.05 (i.e., 5%); instead, amino acids vary in their frequency of occurrence.

Gly	0.089	Arg	0.041
Ala	0.087	Asn	0.040
Leu	0.085	Phe	0.040
Lys	0.081	Gln	0.038
Ser	0.070	Ile	0.037
Val	0.065	His	0.034
Thr	0.058	Cys	0.033
Pro	0.051	Tyr	0.030
Glu	0.050	Met	0.015
Asp	0.047	Trp	0.010

If 20 amino acids occurred in nature at equal frequencies, each would be observed 5% of the time. However, some are more common (G, A, L, K) and some rare (C, Y, M, W).

Normalized frequencies of amino acids:
we need these values to calculate denominator $p_i p_j$

Gly	8.9%	Arg	4.1%
Ala	8.7%	Asn	4.0%
Leu	8.5%	Phe	4.0%
Lys	8.1%	Gln	3.8%
Ser	7.0%	Ile	3.7%
Val	6.5%	His	3.4%
Thr	5.8%	Cys	3.3%
Pro	5.1%	Tyr	3.0%
Glu	5.0%	Met	1.5%
Asp	4.7%	Trp	1.0%

- blue=6 codons; red=1 codon in the genetic code
- These frequencies f_i sum to 1

Dayhoff model step 3: amino acid substitutions

	A Ala	R Arg	N Asn	D Asp	C Cys	Q Gln	E Glu	G Gly	H His	I Ile	L Leu	K Lys	M Met	F Phe	P Pro	S Ser	T Thr	W Trp	Y Tyr	V Val
A																				
R	30																			
N	109	17																		
D	154	0	532																	
C	33	10	0	0																
Q	93	120	50	76	0															
E	266	0	94	831	0	422														
G	579	10	156	162	10	30	112													
H	21	103	226	43	10	243	23	10												
I	66	30	36	13	17	8	35	0	3											
L	95	17	37	0	y	75	15	17	40	253										
K	57	477	322	85	0	147	104	60	23	43	39									
M	29	17	0	0	0	20	7	7	0	57	207	90								
F	20	7	7	0	0	0	0	17	20	90	167	0	17							
P	345	67	27	10	10	93	40	49	50	7	43	43	4	7						
S	772	137	432	98	117	47	86	450	26	20	32	168	20	40	269					
T	590	20	169	57	10	37	31	50	14	129	52	200	28	10	73	696				
W	0	27	3	0	0	0	0	0	3	0	13	0	0	10	0	17	0			
Y	20	3	36	0	30	0	10	0	40	13	23	10	0	260	0	22	23	6		
V	365	20	13	17	33	27	37	97	30	661	303	17	77	10	50	43	186	0	17	
	A Ala	R Arg	N Asn	D Asp	C Cys	Q Gln	E Glu	G Gly	H His	I Ile	L Leu	K Lys	M Met	F Phe	P Pro	S Ser	T Thr	W Trp	Y Tyr	V Val

From a survey of 1572 observed substitutions, the original amino acid (columns) are compared to the changes (rows).

[illegible]

Figure 80. Numbers of accepted point mutations (X10) accumulated from closely related sequences. Fifteen hundred and seventy-

two exchanges are shown. Fractional exchanges result when ancestral sequences are ambiguous.

Dayhoff model step 3: amino acid substitutions

	A Ala	R Arg	N Asn	D Asp	C Cys	Q Gln	E Glu	G Gly
A								
R	30							
N	109	17						
D	154	0	532					
C	33	10	0	0				
Q	93	120	50	76	0			
E	266	0	94	831	0	422		
G	579	10	156	162	10	30	112	
H	21	103	226	43	10	243	23	10

Zooming in on the previous table, note that substitutions are very common (e.g. $D \rightarrow E$, $A \rightarrow G$) while others are rare (e.g. $C \rightarrow Q$, $C \rightarrow E$). The scoring system we use for pairwise alignments should reflect these trends.

Dayhoff model step 3: Relative mutability of amino acids

TABLE 3.2 Relative mutabilities of amino acids. The value of alanine is arbitrarily set to 100.

Asn	134	His	66
Ser	120	Arg	65
Asp	106	Lys	56
Glu	102	Pro	56
Ala	100	Gly	49
Thr	97	Tyr	41
Ile	96	Phe	41
Met	94	Leu	40
Gln	93	Cys	20
Val	74	Trp	18

Dayhoff et al. used (1) data on the frequency of amino acids and (2) data on observed and inferred numbers of substitutions to determine the relative mutability of amino acids. In a scoring system alignment of two tryptophans will be weighted more heavily than two asparagines.

Dayhoff step 4 (of 7): Mutation probability matrix for the evolutionary distance of 1 PAM

		Original amino acid																			
		A Ala	R Arg	N Asn	D Asp	C Cys	Q Gln	E Glu	G Gly	H His	I Ile	L Leu	K Lys	M Met	F Phe	P Pro	S Ser	T Thr	W Trp	Y Tyr	V Val
Replacement amino acid	A	98.7	0.0	0.1	0.1	0.0	0.1	0.2	0.2	0.0	0.1	0.0	0.0	0.1	0.0	0.2	0.4	0.3	0.0	0.0	0.2
	R	0.0	99.1	0.0	0.0	0.0	0.1	0.0	0.0	0.1	0.0	0.0	0.2	0.0	0.0	0.0	0.1	0.0	0.1	0.0	0.0
	N	0.0	0.0	98.2	0.4	0.0	0.0	0.1	0.1	0.2	0.0	0.0	0.1	0.0	0.0	0.0	0.2	0.1	0.0	0.0	0.0
	D	0.1	0.0	0.4	98.6	0.0	0.1	0.5	0.1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.1	0.0	0.0	0.0	0.0
	C	0.0	0.0	0.0	0.0	99.7	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.1	0.0	0.0	0.0	0.0
	Q	0.0	0.1	0.0	0.1	0.0	98.8	0.3	0.0	0.2	0.0	0.0	0.1	0.0	0.0	0.1	0.0	0.0	0.0	0.0	0.0
	E	0.1	0.0	0.1	0.6	0.0	0.4	98.7	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	G	0.2	0.0	0.1	0.1	0.0	0.0	0.1	99.4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.2	0.0	0.0	0.0	0.1
	H	0.0	0.1	0.2	0.0	0.0	0.2	0.0	0.0	99.1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	I	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	98.7	0.1	0.0	0.2	0.1	0.0	0.0	0.1	0.0	0.0	0.3
	L	0.0	0.0	0.0	0.0	0.0	0.1	0.0	0.0	0.0	0.2	99.5	0.0	0.5	0.1	0.0	0.0	0.0	0.0	0.0	0.2
	K	0.0	0.4	0.3	0.1	0.0	0.1	0.1	0.0	0.0	0.0	0.0	99.3	0.2	0.0	0.0	0.1	0.1	0.0	0.0	0.0
	M	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.1	0.1	0.0	98.7	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	F	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.1	0.1	0.0	0.0	99.5	0.0	0.0	0.0	0.0	0.3	0.0
	P	0.1	0.1	0.0	0.0	0.0	0.1	0.0	0.0	0.1	0.0	0.0	0.0	0.0	0.0	99.3	0.1	0.0	0.0	0.0	0.0
	S	0.3	0.1	0.3	0.1	0.1	0.0	0.1	0.2	0.0	0.0	0.0	0.1	0.0	0.0	0.2	98.4	0.4	0.1	0.0	0.0
	T	0.2	0.0	0.1	0.0	0.0	0.0	0.0	0.0	0.0	0.1	0.0	0.1	0.1	0.0	0.1	0.3	98.7	0.0	0.0	0.1
	W	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	99.8	0.0	0.0
	Y	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.2	0.0	0.0	0.0	0.0	99.5	0.0
	V	0.1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.6	0.1	0.0	0.2	0.0	0.0	0.0	0.1	0.0	0.0	99.0

This mutation probability matrix includes original amino acids (columns) and replacements (rows). The diagonals show that at a distance of 1 PAM most residues remain the same about 99% of the time (see shaded entries). Note how cysteine (C) and tryptophan (W) undergo few substitutions, and asparagine (N) many.

Substitution Matrix

- A substitution matrix contains values proportional to the probability that amino acid i mutates into amino acid j for all pairs of amino acids.
- Substitution matrices are constructed by assembling a large and diverse sample of verified pairwise alignments (or multiple sequence alignments) of amino acids.
- Substitution matrices should reflect the true probabilities of mutations occurring through a period of evolution.
- The two major types of substitution matrices are PAM and BLOSUM.

PAM matrices:

Point-accepted mutations

- **PAM matrices** are based on global alignments of closely related proteins.
- **The PAM1 is the matrix calculated from comparisons of sequences with no more than 1% divergence.** At an evolutionary interval of PAM1, one change has occurred over a length of 100 amino acids.
- Other PAM matrices are extrapolated from PAM1. For PAM250, 250 changes have occurred for two proteins over a length of 100 amino acids.
- All the PAM data come from closely related proteins (>85% amino acid identity).

Dayhoff step 4 (of 7): Mutation probability matrix for the evolutionary distance of 1 PAM

		A	R	N	D	C	Q	E	G	H
		Ala	Arg	Asn	Asp	Cys	Gln	Glu	Gly	His
amino acid	A	98.7	0.0	0.1	0.1	0.0	0.1	0.2	0.2	0.0
	R	0.0	99.1	0.0	0.0	0.0	0.1	0.0	0.0	0.1
	N	0.0	0.0	98.2	0.4	0.0	0.0	0.1	0.1	0.2
	D	0.1	0.0	0.4	98.6	0.0	0.1	0.5	0.1	0.0
	C	0.0	0.0	0.0	0.0	99.7	0.0	0.0	0.0	0.0
	Q	0.0	0.1	0.0	0.1	0.0	98.8	0.3	0.0	0.2
	E	0.1	0.0	0.1	0.6	0.0	0.4	98.7	0.0	0.0
	G	0.2	0.0	0.1	0.1	0.0	0.0	0.1	99.4	0.0
	H	0.0	0.1	0.2	0.0	0.0	0.2	0.0	0.0	99.1

At this evolutionary distance of 1 PAM, 1% of the amino acids have diverged between each pair of sequences. The columns are percentages that sum to 100%.

Dayhoff step 5 (of 7): PAM250 and other PAM matrices

<u>NP 002037.2</u>	164	IHDNFGIVEGLMTTVHAIITATQKTVDGPGSGKLWRDGRGALQNII	207
<u>XP 001162057.1</u>	164	IHDNFGIVEGLMTTVHAIITATQKTVDGPGSGKLWRDGRGALQNII	207
<u>NP 001003142.1</u>	162	IHDHFGIVEGLMTTVHAIITATQKTVDGPGSGKMWRDGRGAAQNII	205
<u>XP 893121.1</u>	168	IHDNFGIMEGLMTTVHAIITATQKTVDGPGSGKLWRDGRGAAQNII	211
<u>XP 576394.1</u>	162	IHDNFGIVEGLMTTVHAIITATQKTVDGPGSGKLWRDGRGAAQNII	205
<u>NP 058704.1</u>	162	IHDNFGIVEGLMTTVHAIITATQKTVDGPGSGKLWRDGRGAAQNII	205
<u>XP 001070653.1</u>	162	IHDNFGIVEGLMTTVHAIITATQKTVDGPGSGKLWRDGRGAAQNII	205
<u>XP 001062726.1</u>	162	IHDNFGIVEGLMTTVHAIITATQKTVDGPGSGKLWRDGRGAAQNII	205
<u>NP 989636.1</u>	162	IHDNFGIVEGLMTTVHAIITATQKTVDGPGSGKLWRDGRGAAQNII	205
<u>NP 525091.1</u>	161	INDNFEIVEGLMTTVHATTATQKTVDGPGSGKLWRDGRGAAQNII	204
<u>XP 318655.2</u>	161	INDNFGILEGLMTTVHATTATQKTVDGPGSGKLWRDGRGAAQNII	204
<u>NP 508535.1</u>	170	INDNFGIIEGLMTTVHAVTATQKTVDGPGSGKLWRDGRGAGQNII	213
<u>NP 595236.1</u>	164	INDTFGIEEGLMTTVHATTATQKTVDGPGSKDWRGGRGASANII	207
<u>NP 011708.1</u>	162	INDAFGIEEGLMTTVHSLTATQKTVDGPGSHKDWRGGRTASGNII	205
<u>XP 456022.1</u>	161	INDEFGIDEALMTTVHSITATQKTVDGPGSHKDWRGGRTASGNII	204
<u>NP 001060897.1</u>	166	IHDNFGIIEGLMTTVHAIITATQKTVDGPGSSKDWRGGRAASFNII	209

Consider a multiple alignment of glyceraldehyde 3-phosphate protein sequences. Some substitutions are observed in columns (arrowheads). These give us insight into changes tolerated by natural selection.

Dayhoff step 5 (of 7): PAM250 and other PAM matrices

	▼		▼	▼		▼	▼	▼		▼																																																
mouse	A	I	P	N	P	S	F	L	A	M	P	T	N	E	N	Q	D	N	T	A	I	P	T	I	D	P	I	T	P	I	V	S	T	--	P	V	P	T	M	-----	E	S	I	V	N	T	V	A	N	P	E	A	S	T				
rabbit	S	--	H	P	F	F	M	A	I	L	P	N	K	M	Q	D	K	A	V	T	P	T	T	N	T	I	A	A	V	E	P	T	--	P	I	P	T	T	-----	E	P	V	V	S	T	E	V	I	A	E	A	S	P					
sheep	P	H	P	H	L	S	F	M	A	I	P	P	K	K	D	Q	D	K	T	E	I	P	A	I	N	T	I	A	S	A	E	P	T	V	H	S	T	P	T	-----	E	A	V	V	N	A	V	D	N	P	E	A	S	S				
cattle	P	H	P	H	L	S	F	M	A	I	P	P	K	K	N	Q	D	K	T	E	I	P	T	I	N	T	I	A	S	G	E	P	T	--	S	T	P	T	-----	E	A	V	E	S	T	V	A	T	L	E	D	S	P					
pig	P	R	P	H	A	S	F	I	A	I	P	P	K	K	N	Q	D	K	T	A	I	P	A	I	N	S	I	A	T	V	E	P	T	--	I	V	P	A	T	E	P	I	V	N	A	E	P	I	V	N	A	V	T	P	E	A	S	S
human	P	N	L	H	P	S	F	I	A	I	P	P	K	K	I	Q	D	K	I	I	I	P	T	I	N	T	I	A	T	V	E	P	T	--	P	A	P	A	T	-----	E	P	T	V	D	S	V	V	T	P	E	A	F	S				
horse	P	C	P	H	P	S	F	I	A	I	P	P	K	K	L	Q	E	I	T	V	I	P	K	I	N	T	I	A	T	V	E	P	T	--	P	I	P	T	P	-----	E	P	T	V	N	N	A	V	I	P	D	A	S	S				
	.	:	*	:	*	:	.	:	:	*	:	*	:	.	:	*	:	.	*	:	.	*	:	.	*	:	.	*	:	.	*	:	.	*	:	.	*	:	.	*	:	.	*	:	.	*	:	.	*	:	.							

Now consider the alignment of distantly related kappa caseins. There are few conserved column positions, and many some columns (double arrowheads) have five different residues among the 7 proteins.

We want to design a scoring system that is tolerant of distantly related proteins: if the scoring system is too strict then the divergent sequences may be penalized so heavily that authentic homologs are not identified or aligned.

Dayhoff step 5 (of 7): PAM250 and other PAM matrices

replacement amino acid	original amino acid								
	PAM0	A	R	N	D	C	Q	E	G
	A	100	0	0	0	0	0	0	0
	R	0	100	0	0	0	0	0	0
	N	0	0	100	0	0	0	0	0
	D	0	0	0	100	0	0	0	0
	C	0	0	0	0	100	0	0	0
	Q	0	0	0	0	0	100	0	0
	E	0	0	0	0	0	0	100	0
	G	0	0	0	0	0	0	0	100
	original amino acid								
	PAM ∞	A	R	N	D	C	Q	E	G
	A	8.7	8.7	8.7	8.7	8.7	8.7	8.7	8.7
	R	4.1	4.1	4.1	4.1	4.1	4.1	4.1	4.1
	N	4.0	4.0	4.0	4.0	4.0	4.0	4.0	4.0
	D	4.7	4.7	4.7	4.7	4.7	4.7	4.7	4.7
	C	3.3	3.3	3.3	3.3	3.3	3.3	3.3	3.3
	Q	3.8	3.8	3.8	3.8	3.8	3.8	3.8	3.8
	E	5.0	5.0	5.0	5.0	5.0	5.0	5.0	5.0
	G	8.9	8.9	8.9	8.9	8.9	8.9	8.9	8.9

At the extreme of perfectly conserved proteins (PAM0) there are no amino acid replacements. At the extreme of completely diverged proteins (PAM ∞) the matrix converges on the background frequencies of the amino acids.

PAM250 matrix: for proteins that share ~20% identity

		Original amino acid																			
		A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
Replacement amino acid	A	13	6	9	9	5	8	9	12	6	8	6	7	7	4	11	11	11	2	4	9
	R	3	17	4	3	2	5	3	2	6	3	2	9	4	1	4	4	3	7	2	2
	N	4	4	6	7	2	5	6	4	6	3	2	5	3	2	4	5	4	2	3	3
	D	5	4	8	11	1	7	10	5	6	3	2	5	3	1	4	5	5	1	2	3
	C	2	1	1	1	52	1	1	2	2	2	1	1	1	1	2	3	2	1	4	2
	Q	3	5	5	6	1	10	7	3	7	2	3	5	3	1	4	3	3	1	2	3
	E	5	4	7	11	1	9	12	5	6	3	2	5	3	1	4	5	5	1	2	3
	G	12	5	10	10	4	7	9	27	5	5	4	6	5	3	8	11	9	2	3	7
	H	2	5	5	4	2	7	4	2	15	2	2	3	2	2	3	3	2	2	3	2
	I	3	2	2	2	2	2	2	2	2	10	6	2	6	5	2	3	4	1	3	9
	L	6	4	4	3	2	6	4	3	5	15	34	4	20	13	5	4	6	6	7	13
	K	6	18	10	8	2	10	8	5	8	5	4	24	9	2	6	8	8	4	3	5
	M	1	1	1	1	0	1	1	1	1	2	3	2	6	2	1	1	1	1	1	2
	F	2	1	2	1	1	1	1	1	3	5	6	1	4	32	1	2	2	4	20	3
	P	7	5	5	4	3	5	4	5	5	3	3	4	3	2	20	6	5	1	2	4
	S	9	6	8	7	7	6	7	9	6	5	4	7	5	3	9	10	9	4	4	6
	T	8	5	6	6	4	5	5	6	4	6	4	6	5	3	6	8	11	2	3	6
	W	0	2	0	0	0	0	0	0	1	0	1	0	0	1	0	1	0	55	1	0
	Y	1	1	2	1	3	1	1	1	3	2	2	1	2	15	1	2	2	3	31	2
	V	7	4	4	4	4	4	4	5	4	15	10	4	10	5	5	5	7	2	4	17

Compare this to a PAM1 matrix, and note the diagonal still has high scores but much information content is lost.

Table 23
Correspondence between Observed Differences
and the Evolutionary Distance

Observed Percent Difference	Evolutionary Distance in PAMs
1	1
5	5
10	11
15	17
20	23
25	30
30	38
35	47
40	56
45	67
50	80
55	94
60	112
65	133
70	159
75	195
80	246
85	328

Dayhoff step 6 (of 7): from a mutation probability matrix to a relatedness odds matrix

$$R_{ij} = \frac{M_{ij}}{f_i}$$

1. A **relatedness odds matrix** reports the probability that amino acid *j* will change to *i* in a homologous sequence.
2. The numerator models the observed change. The denominator *f_i* is the probability of amino acid residue *i* occurring in the second sequence by chance.
3. A positive value indicates a replacement happens more often than expected by chance. A **negative value** indicates the **replacement is not favored**.

Why do we go from a mutation probability matrix to a log odds matrix?

- We want a scoring matrix so that when we do a pairwise alignment (or a BLAST search) we know what score to assign to two aligned amino acid residues.
- Logarithms are easier to use for a scoring system. They allow us to sum the scores of aligned residues (rather than having to multiply them).

How do we go from a mutation probability matrix to a log odds matrix?

The cells in a log odds matrix consist of an “odds ratio”:

the probability that an alignment is authentic
the probability that the alignment was random

The score S for an alignment of residues a,b is given by:

$$S(a,b) = 10 \log_{10} (M_{ab}/p_b)$$

As an example, for tryptophan,

$$s_{i,j} = 10 \times \log \left(\frac{q_{i,j}}{p_{i,j}} \right)$$

$$S(\text{trp}, \text{trp}) = 10 \log_{10} (0.55/0.010) = 17.4$$

What do the numbers mean in a log odds matrix?

A score of **+2 indicates** that the amino acid replacement occurs **1.6 times as frequently as expected by chance**.

A score of **0 is neutral**.

A score of **-10 indicates** that the correspondence of two amino acids in an alignment that accurately represents homology (evolutionary descent) **is one tenth as frequent as the chance alignment** of these amino acids.

Dayhoff step 7 (of 7): log odds scoring matrix

$$s_{ij} = 10 \times \log_{10} \left(\frac{M_{ij}}{f_i} \right).$$

- A log odds matrix is the logarithmic form of the relatedness odds matrix.
- s_{ij} is the score for aligning any two residues in a pairwise alignment. (There is also a score for aligning a residue with itself.)
- M_{ij} is of the observed frequency of substitutions for each pair of amino acids. These values (“target frequencies”) are derived from a mutation probability matrix.

Example of a score for aligning cysteine and leucine using the values in a PAM250 scoring matrix

$$s_{(\text{cysteine, leucine})} = 10 \times \log_{10} \left(\frac{0.02}{0.085} \right) = -6.3$$

Log-odds matrix for PAM250

6									
-3	5								
4	0	6							
2	-5	0	9						
-3	-1	-2	-5	6					
-3	0	-2	-3	1	2				
-2	0	-1	-3	0	1	3			
-2	-3	-4	0	-6	-2	-5	17		
-1	-4	-2	7	-5	-3	-3	0	10	
2	-2	2	-1	-1	-1	0	-6	-2	4
L	K	M	F	P	S	T	W	Y	V

This is a useful matrix for comparing **distantly related proteins**. Note that an alignment of two tryptophan (W) residues earns +17 and a W to T mismatch is -5.

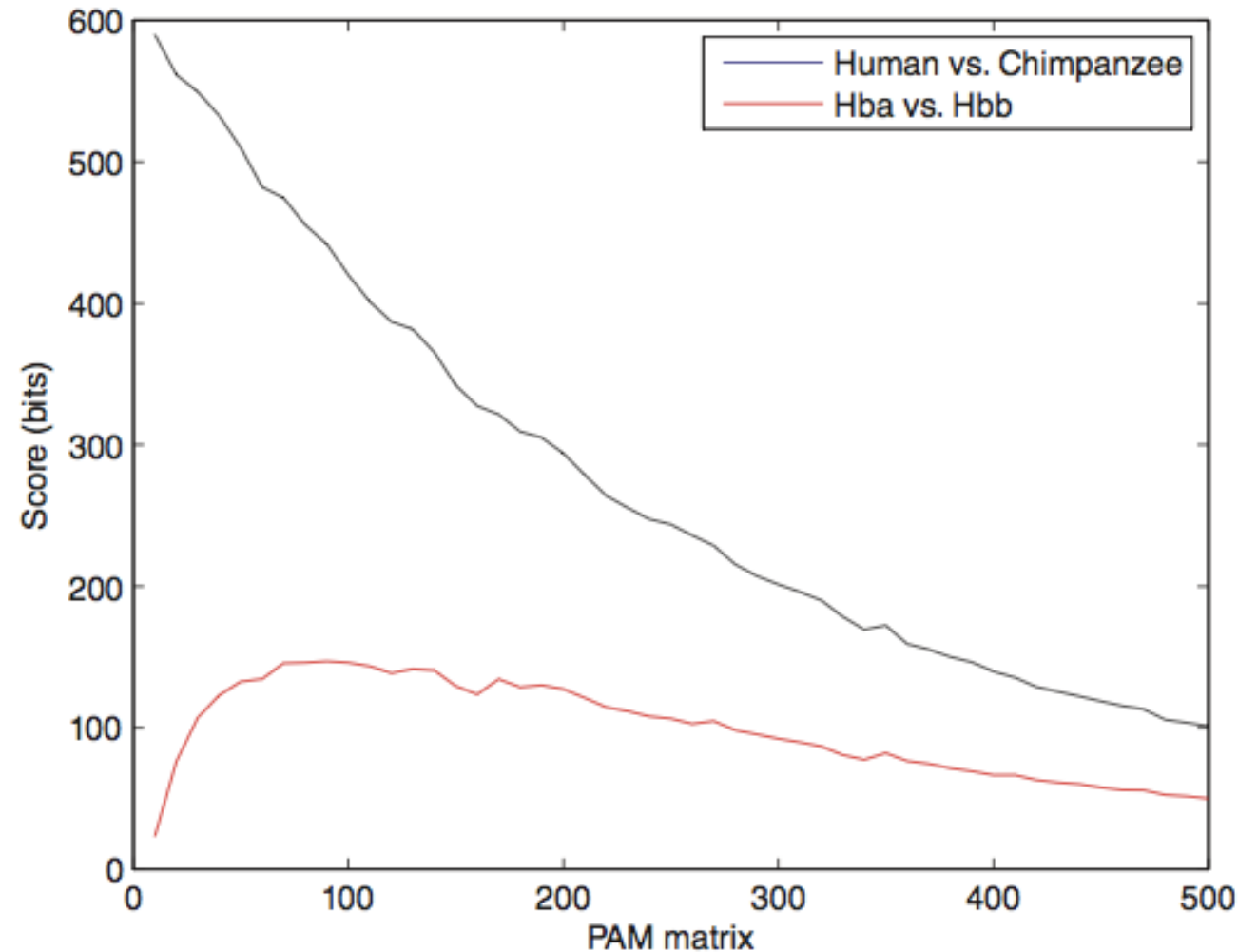
Log-odds matrix for PAM10

A	7																			
R	-10	9																		
N	-7	-9	9																	
D	-6	-17	-1	8																
C	-10	-11	-17	-21	10															
Q	-7	-4	-7	-6	-20	9														
E	-5	-15	-5	0	-20	-1	8													
G	-4	-13	-6	-6	-13	-10	-7	7												
H	-11	-4	-2	-7	-10	-2	-9	-13	10											
I	-8	-8	-8	-11	-9	-11	-8	-17	-13	9										
L	-9	-12	-10	-19	-21	-8	-13	-14	-9	-4	7									
K	-10	-2	-4	-8	-20	-6	-7	-10	-10	-9	-11	7								
M	-8	-7	-15	-17	-20	-7	-10	-12	-17	-3	-2	-4	12							
F	-12	-12	-12	-21	-19	-19	-20	-12	-9	-5	-5	-20	-7	9						
P	-4	-7	-9	-12	-11	-6	-9	-10	-7	-12	-10	-10	-11	-13	8					
S	-3	-6	-2	-7	-6	-8	-7	-4	-9	-10	-12	-7	-8	-9	-4	7				
T	-3	-10	-5	-8	-11	-9	-9	-10	-11	-5	-10	-6	-7	-12	-7	-2	8			
W	-2	-5	-11	-21	-22	-19	-23	-21	-10	-20	-9	-18	-19	-7	-20	-8	-19	13		
Y	-11	-14	-7	-17	-7	-18	-11	-20	-6	-9	-10	-12	-17	-1	-20	-10	-9	-8	10	
V	-5	-11	-12	-11	-9	-10	-10	-9	-9	-1	-5	-13	-4	-12	-9	-10	-6	-22	-10	8
	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V

More closely related proteins

This is an example of a scoring matrix **with “severe” penalties**. A match of W to W earns +13, but a mismatch (e.g. W aligned to T) has a score of -19, far lower than in PAM250.

Effect of scoring matrix on bit scores



Look at score for distantly related proteins (e.g. beta globin versus alpha globin) and note that PAM10 or similar matrices assign very low scores. This effect is not seen for very closely related proteins (e.g. a chimp vs. human globin).

PAM matrices are based on data from the alignment of closely related protein families, and they involve the assumption that substitution probabilities for highly related proteins (e.g., PAM40) can be extrapolated to probabilities for distantly related proteins (e.g., PAM250).

BLOSUM matrices are based on empirical observations of more distantly related protein alignments.

https://www.youtube.com/watch?v=0_66UK-439M&t=13s

BLOSUM62 scoring matrix

blocks substitution matrix

A	4																			
R	-1	5																		
N	-2	0	6																	
D	-2	-2	1	6																
C	0	-3	-3	-3	9															
Q	-1	1	0	0	-3	5														
E	-1	0	0	2	-4	2	5													
G	0	-2	0	-1	-3	-2	-2	6												
H	-2	0	1	-1	-3	0	0	-2	8											
I	-1	-3	-3	-3	-1	-3	-3	-4	-3	4										
L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4									
K	-1	2	0	-1	-1	1	1	-2	-1	-3	-2	5								
M	-1	-2	-2	-3	-1	0	-2	-3	-2	1	2	-1	5							
F	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6						
P	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7					
S	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4				
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5			
W	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11		
Y	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7	
V	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4
	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V

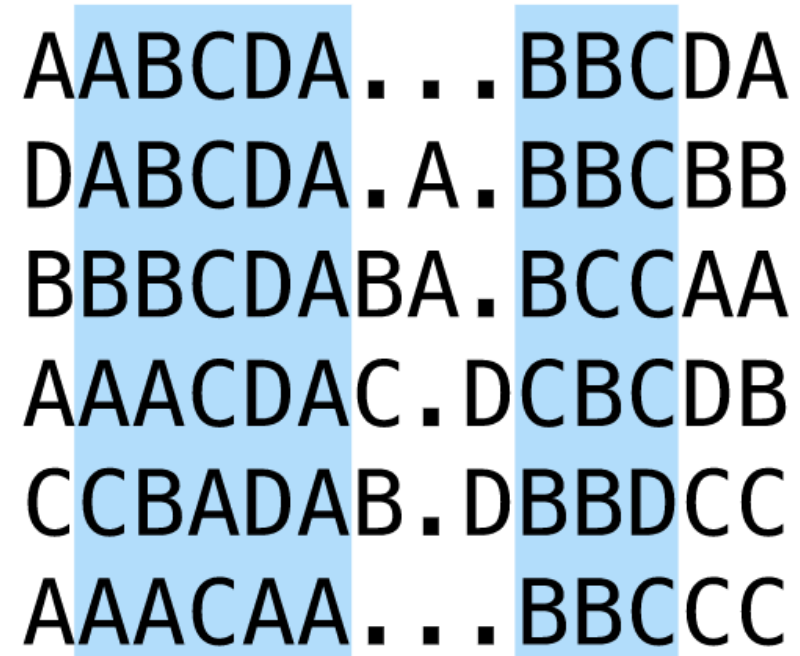
BL62 is the default scoring matrix at the NCBI BLAST site.

Derivation of BLOSUM matrices

BLOSUM matrices are derived from comparisons of blocks of sequences from the Blocks database.

What are blocks and what is the blocks database?

A *block* is an ungapped multiple alignments of highly conserved, short regions. Here is what a sample block looks like:



```
AABCDAA . . . BBCDA
DABCDAA . A . BBCBB
BBBCDABA . BCCAA
AAACDAC . DCBCDB
CCBADAB . DBBDC
AAACAA . . . BBCCC
```

Conserved blocks in alignment

The blocks database contains multiple alignments of conserved regions in protein families.

<https://snipcademy.com/pairwise-alignment#blosum---blocks-substitution-matrix>

BLOSUM Matrices

BLOSUM matrices are based on local alignments.

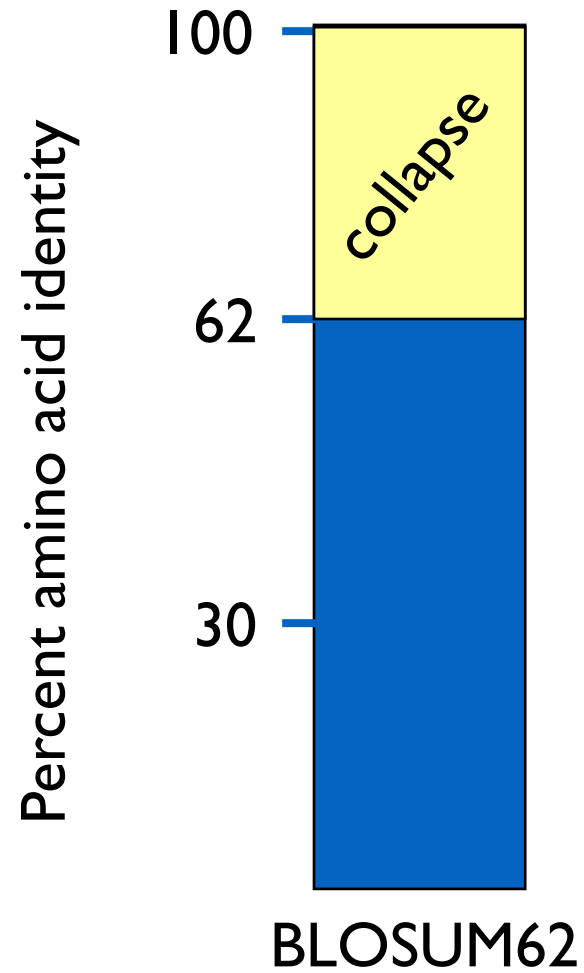
All BLOSUM matrices are based on observed alignments; they are not extrapolated from comparisons of closely related proteins.

BLOSUM stands for blocks substitution matrix.

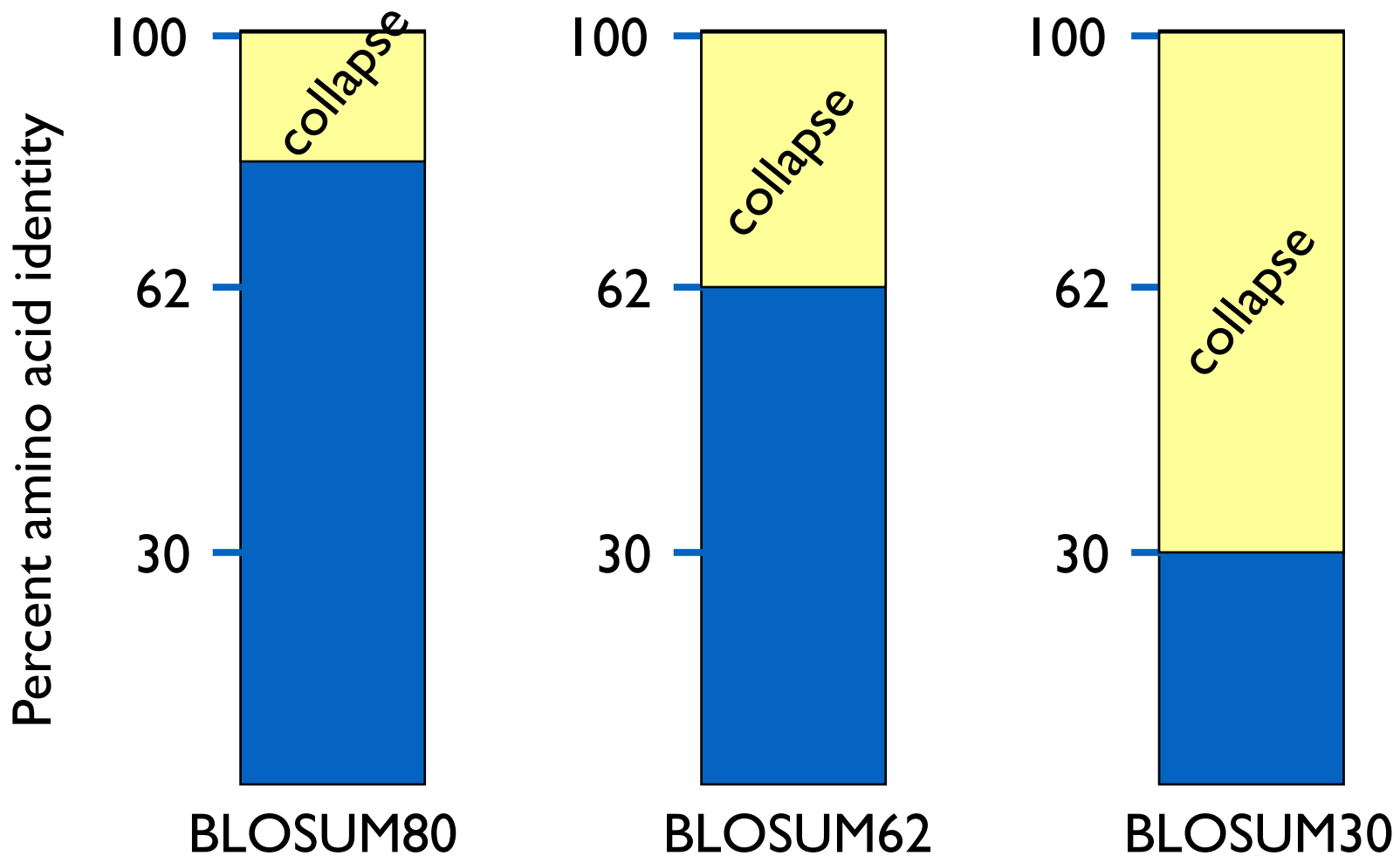
BLOSUM62 is a matrix calculated from comparisons of sequences with no less than 62% divergence.

BLOSUM62 is the default matrix in BLAST 2.0.

BLOSUM Matrices



BLOSUM Matrices



PAM matrices with similar BLOSUM matrices (relative entropy of each PAM matrix is from Altschul 1991)^[18]

PAM matrix	Equivalent BLOSUM matrix	Relative entropy (bits)
PAM100	Blosum90	1.18
PAM120	Blosum89	0.98
PAM160	Blosum60	0.70
PAM200	Blosum52	0.51
PAM250	Blosum45	0.36

Source: Wikipedia

Summary of PAM and BLOSUM matrices

BLOSUM90

BLOSUM62

BLOSUM45

PAM30

PAM120

PAM250

Less divergent



More divergent

Human versus
chimpanzee beta globin

Human versus
bacterial globins

A higher PAM number, and a lower BLOSUM number, tends to correspond to a matrix tuned to more divergent proteins.