

Substitution matrices

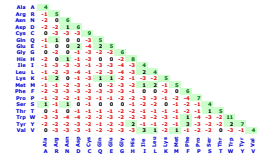
Jacques van Helden

Jacques.van-Helden@univ-amu.fr
 Université d'Aix-Marseille, France
 Lab. Technological Advances for Genomics and Clinics
 (TAGC, INSERM Unit U1090)
<http://tagc.univ-mra.fr/>

FORMER ADDRESS (1999-2011)
 Université Libre de Bruxelles, Belgique
 Bioinformatic des Génomes et des Réseaux (BIGRe lab)
<http://www.bigre.ulb.ac.be/>



	A	C	G	T
A	2			
C	-2	2		
G	-2	-2	2	
T	-1	-2	-2	2



- A **substitution matrix** indicates the score associated to each possible pair of aligned residues in an alignment.
 - Each row and each column represent one of the possible residues (4 for DNA, 20 for proteins).
 - The diagonal indicates identities.
 - The lower triangle indicates substitutions.
 - The upper triangle is symmetrical to the lower triangle, and does not need to be displayed.
 - Positive scores indicate that the aligned pair of residue is considered "beneficial" for the alignment. Note that some mismatches might have positive scores in protein alignments (see later).
 - Negative scores are considered as penalties associated to mismatches, and alignment algorithms will try to avoid aligning the pairs of residues.
- One could decide to give a lower cost to A-T substitutions, if we assume that these are more likely to occur in our sequences
- Example: the top matrix represents arbitrarily defined scores for DNA alignment
 - match 2
 - A-T mismatch -1
 - other mismatch -2
- The scoring scheme can be represented as a substitution matrix

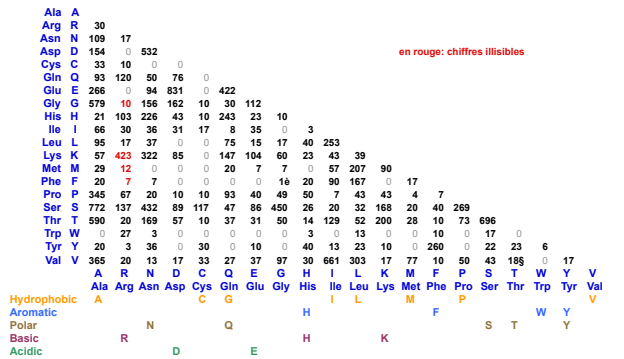
Scoring an alignment matrix with a substitution matrix

- Let us come back to our previous alignment matrix
- For each cell of the alignment matrix, we compare the residue in sequences A and B, and take the score for this pair of residues in the substitution matrix.

	A	C	G	T
A	2			
C	-2	2		
G	-2	-2	2	
T	-1	-2	-2	2

	A	A	T	C	T	T	C	A	G	C	G	T	A	T	T	G	C	T
A	2	-1	-2	-1	-1	-2	-2	-2	-2	-1	2	-1	-1	-2	-2	-1		
T	-1	-1	2	2	2	-1	-2	-2	-2	-2	-1	2	2	2	-2	-2		
C	-2	-2	-2	2	-2	-2	2	2	2	-2	-2	-2	-2	-2	-2	2		
T	-1	-1	2	-2	2	2	-2	-2	-2	-2	-1	2	2	2	-2	-2		
A	2	2	-1	-2	-1	-1	-2	2	-2	-2	-1	2	-1	-1	-2	-2		
G	-2	-2	-2	-2	-2	-2	2	2	2	-2	-2	-2	-2	-2	-2	2		
C	-2	-2	-2	2	-2	-2	2	2	2	-2	-2	-2	-2	-2	-2	2		
G	-2	-2	-2	-2	-2	-2	2	2	2	-2	-2	-2	-2	-2	-2	2		
B	-2	-2	-2	-2	-2	-2	2	2	2	-2	-2	-2	-2	-2	-2	2		
A	2	2	-1	-2	-1	-1	-2	2	-2	-2	-1	2	-1	-1	-2	-2		
G	-2	-2	-2	-2	-2	-2	2	2	2	-2	-2	-2	-2	-2	-2	2		
G	-2	-2	-2	-2	-2	-2	2	2	2	-2	-2	-2	-2	-2	-2	2		
T	-1	-1	2	-2	2	2	-2	-2	-2	-2	-1	2	2	2	-2	-2		
A	2	2	-1	-2	-1	-1	-2	2	-2	-2	-1	2	-1	-1	-2	-2		
T	-1	-1	2	-2	2	2	-2	-2	-2	-2	-1	2	2	2	-2	-2		

Substitution counts in 71 groups of aligned proteins (Dayhoff, 1978)



Substitution matrices for proteins

$$S_{i,j} = s_{j,i} = \log_2 \left(\frac{f_{i,j}}{(f_i f_j)} \right)$$

	C	S	T	P	A	G	...
C	11.5						...
S	0.1	2.2					...
T	-0.5	1.5	2.5				...
P	-3.1	0.4	0.1	7.6			...
A	0.5	1.1	0.6	0.3	2.4		...
G	-2.0	0.4	-1.1	-1.6	0.5	1.6	...
...

- Margaret Dayhoff (1978) measured the rate of substitutions between each pair of amino acids, in a collection of aligned proteins.
- Scores are calculated as **log-odds**
 - Positive values** reflect frequent ("accepted") substitutions, i.e. substitutions that occur more frequently than expected by chance.
 - Negative values** reflect rare ("unfavourable") mutations, i.e. substitutions that occur less frequently than expected by chance
- The diagonal reflect residue conservation

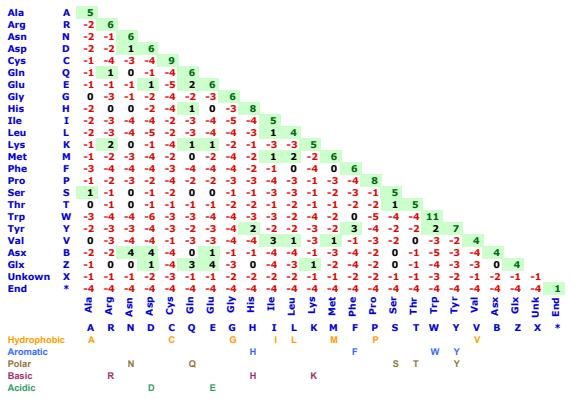
PAM scoring matrices

- The alignments used by Dayhoff had ~85% identity
- However, frequencies of substitutions are expected to depend on the rate of divergence between sequences: the number of substitutions increases with time.
- In order to take into account the divergence rate, Margaret Dayhoff calculated a series of **scoring matrices**, each reflecting a certain level of divergence
 - PAM001 rates of substitutions between amino-acid pairs expected for proteins with an average of 1% substitution per position
 - PAM050 rates of substitutions between amino-acid pairs expected for proteins with an average of 50% substitution per position
 - PAM250 250% mutations/position (note: a position could mutate several times)
- The substitution matrix must this be chosen according to the relatedness of the sequences to be aligned

Reference: Dayhoff et al. (1978). A model of evolutionary change in proteins. In *Atlas of Protein Sequence and Structure*, vol. 5, suppl. 3, 345-352. National Biomedical Research Foundation, Silver Spring, MD, 1978.

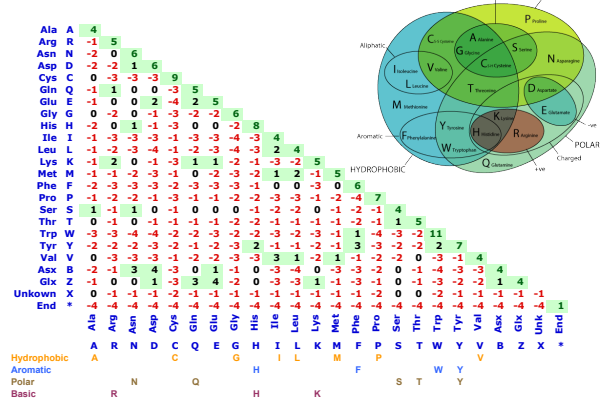
Reference: Dayhoff et al. (1978). A model of evolutionary change in proteins. In *Atlas of Protein Sequence and Structure*, vol. 5, suppl. 3, 345-352. National Biomedical Research Foundation, Silver Spring, MD, 1978.

BLOSUM80



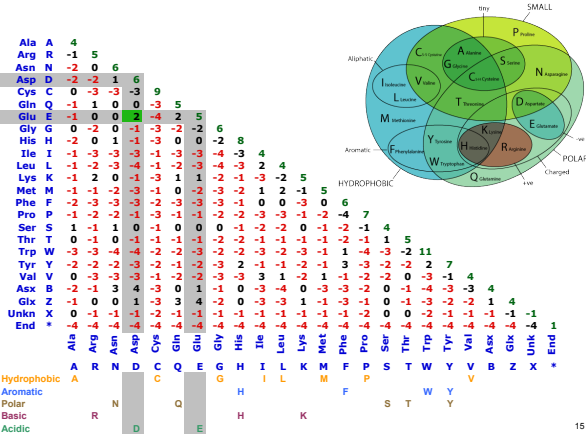
13

BLOSUM62 – Amino acid properties



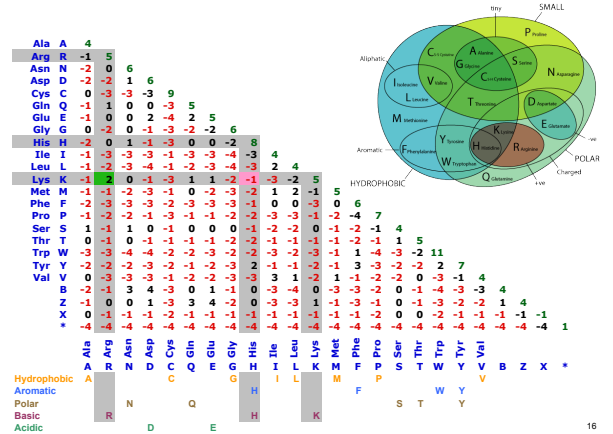
14

BLOSUM62 - substitutions between acidic residues



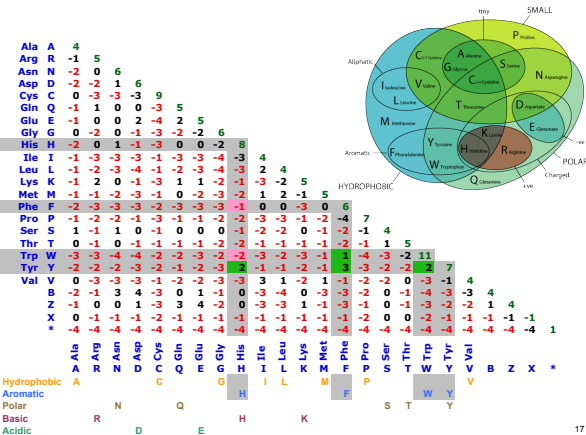
15

BLOSUM62 - substitutions between basic residues



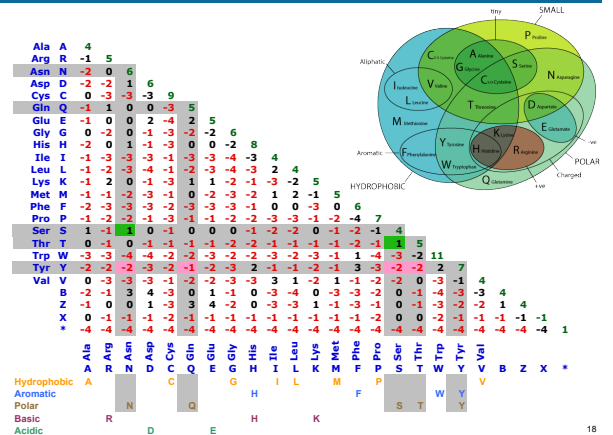
16

BLOSUM62 - substitutions between aromatic residues



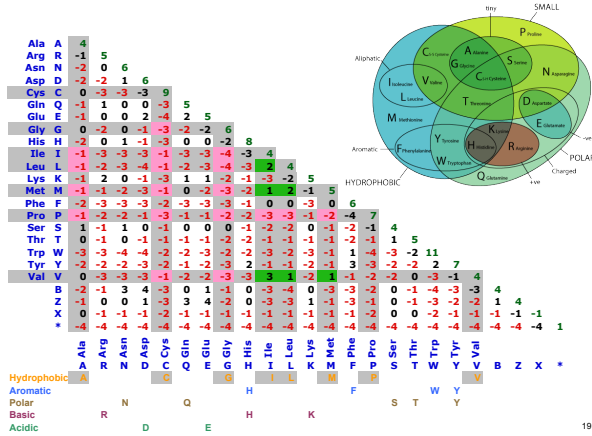
17

BLOSUM62 - substitutions between polar residues



18

BLOSUM62 - substitutions between hydrophobic residues



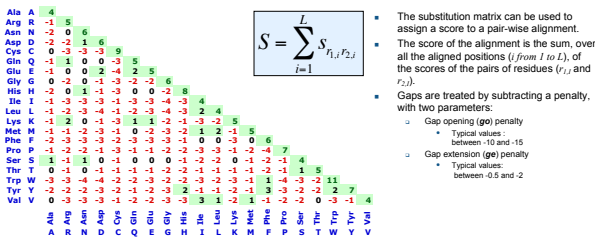
19

Substitution matrices - summary

- Different substitution scoring matrices have been established
 - Residue categories (Phylip)
 - PAM (Dayhoff, 1979).
 - PAM means "Percent Accepted Mutations"
 - BLOSUM (Henikoff & Henikoff, 1992).
 - BLOSUM means "Block sum".
- Substitution matrices allow to detect similarities between more distant proteins than what would be detected with the simple identity of residues.
- The matrix must be chosen carefully, depending on the expected rate of conservation between the sequences to be aligned.
- Beware
 - With PAM matrices
 - the score indicates the percentage of substitution per position
 - > **higher numbers** are appropriate for **more distant** proteins
 - With BLOSUM matrices
 - the score indicates the percentage of conservation
 - > **higher numbers** are appropriate for **more conserved** proteins

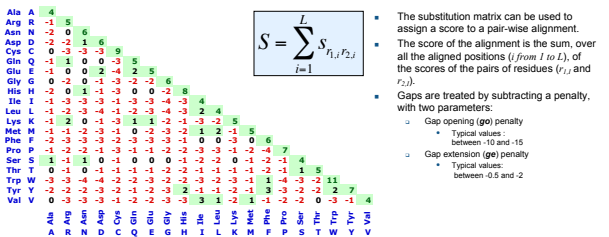
20

Scoring an alignment with a substitution matrix



21

Scoring an alignment with a substitution matrix



22

Bibliography

- Substitution matrices
 - PAM series
 - Dayhoff, M. O., Schwartz, R. M. & Orcutt, B. (1978). A model of evolutionary change in proteins. *Atlas of Protein Sequence and Structure* 5, 345-352.
 - BLOSUM substitution matrices
 - Henikoff, S. & Henikoff, J. G. (1992). Amino acid substitution matrices from protein blocks. *Proc Natl Acad Sci U S A* 89, 10915-9.
 - Gonnet matrices, built by an iterative procedure
 - Gonnet, G. H., Cohen, M. A. & Benner, S. A. (1992). Exhaustive matching of the entire protein sequence database. *Science* 256, 1443-5. 1.

23