

Introduction à la bioinformatique

Matrices de substitution

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://jacques.van-helden.perso.luminy.univ-amu.fr/>

Mismatches et substitutions

- Quand on aligne deux ou plusieurs séquences, on observe souvent des résidus différents à la même position de l'alignement (« **mismatches** »), qui reflètent vraisemblablement qu'une substitution est survenue au sein de l'une des séquences ancestrales.
- On constate que certaines substitutions sont plus fréquentes que d'autres.
- Dans les séquences protéiques, les substitutions fréquentes correspondent généralement à des acides aminés qui partagent des propriétés chimiques (hydrophobie, polarité) ou stérique (encombrement du radical).
- Sur base de cette observation, on construit des **matrices de substitutions** qui serviront ensuite à pondérer les « mismatches » lors de l'alignement de nouvelles séquences.

Construction des matrices de substitutions

***La série de matrices PAM
(« point accepted mutation »)
construite à partir d'alignements par paire***

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) a mesuré les taux de substitutions entre chaque paire d'acides aminés, dans une collection de 71 alignements de paires de protéines.
- A partir des comptages bruts, elle dérive un score de **log-odds**
 - f_i, f_j : fréquences des résidus i et j , resp.
 - $f_{i,j}$: fréquence de la substitution $i \leftrightarrow j$
 - **Les valeurs positives** indiquent des substitutions fréquentes ("acceptées"), c'est-à-dire des substitutions observées plus fréquemment que ce à quoi l'on s'attendrait par hasard.
 - **Les valeurs négatives** indiquent les mutations rares, c'est-à-dire celles qu'on observe moins fréquemment que ce à quoi l'on s'attendrait par hasard. Ce taux inférieur est interprété comme un indice de contre-sélection, suggérant que ces mutations sont généralement défavorables pour la fonction de la protéine.
- La diagonale reflète le taux de conservation des résidus. Notons que certains résidus rares ont un score de conservation très important: le score de conservation n'est pas proportionnel à la fréquence.

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.

PAM scoring matrices

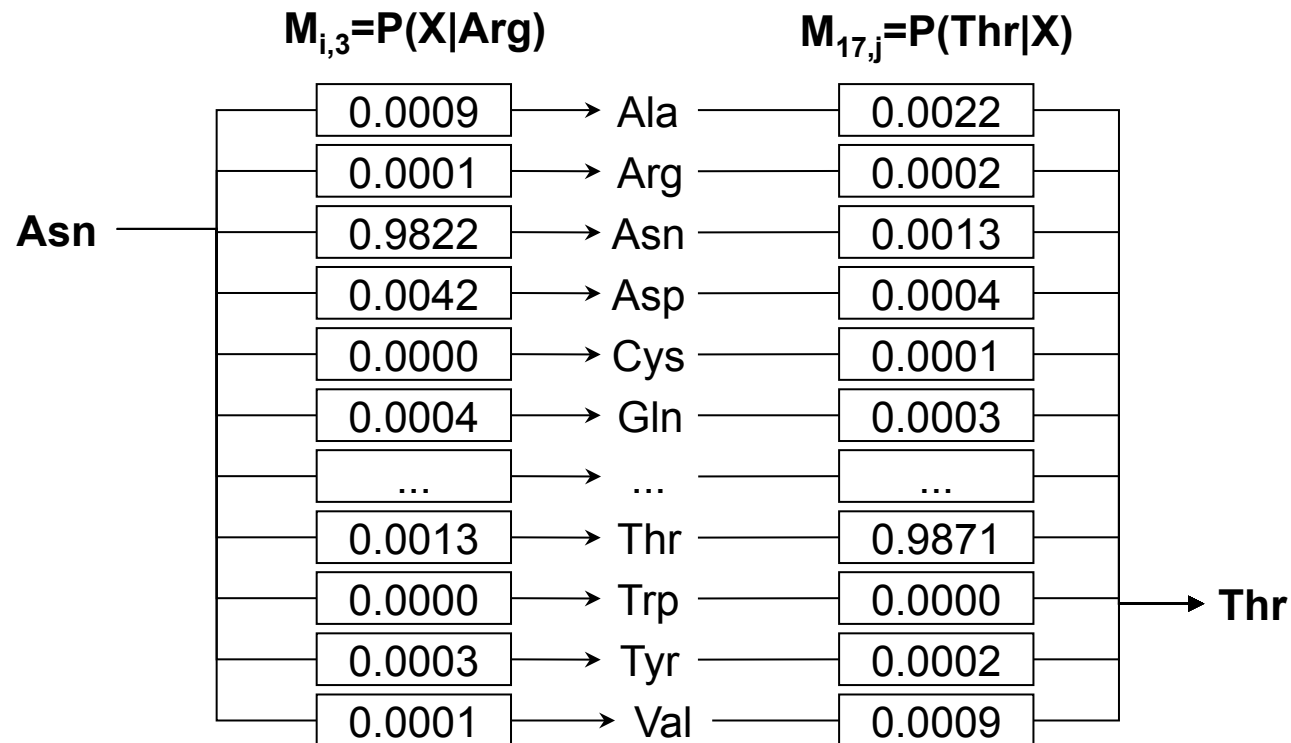
- Les alignements réalisés par Margret Dayhoff en 1987 comportaient un taux moyen d'identité de ~85%.
- Cependant, on s'attend à ce que les fréquences de substitutions dépendent du degré de divergence entre séquences, car leur nombre augmente avec le temps.
- Pour prendre en compte le taux de divergence, Margret Dayhoff a calculé une **série de matrices de score**, reflétant chacune un certain taux de substitutions.

PAM001	taux de substitutions entre acides aminés au terme d'un temps évolutif donnant lieu à ~1% de substitutions par position.
PAM050	taux de substitutions entre acides aminés au terme d'un temps évolutif donnant lieu à ~50% de substitutions par position.
PAM250	idem avec 250% mutations/position (note : une même position peut faire l'objet de plusieurs mutations successives)
- Quand on fait un alignement, on doit choisir l'une des matrices de cette série, en tenant compte du taux de différences entre les deux séquences qu'on veut aligner.

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.

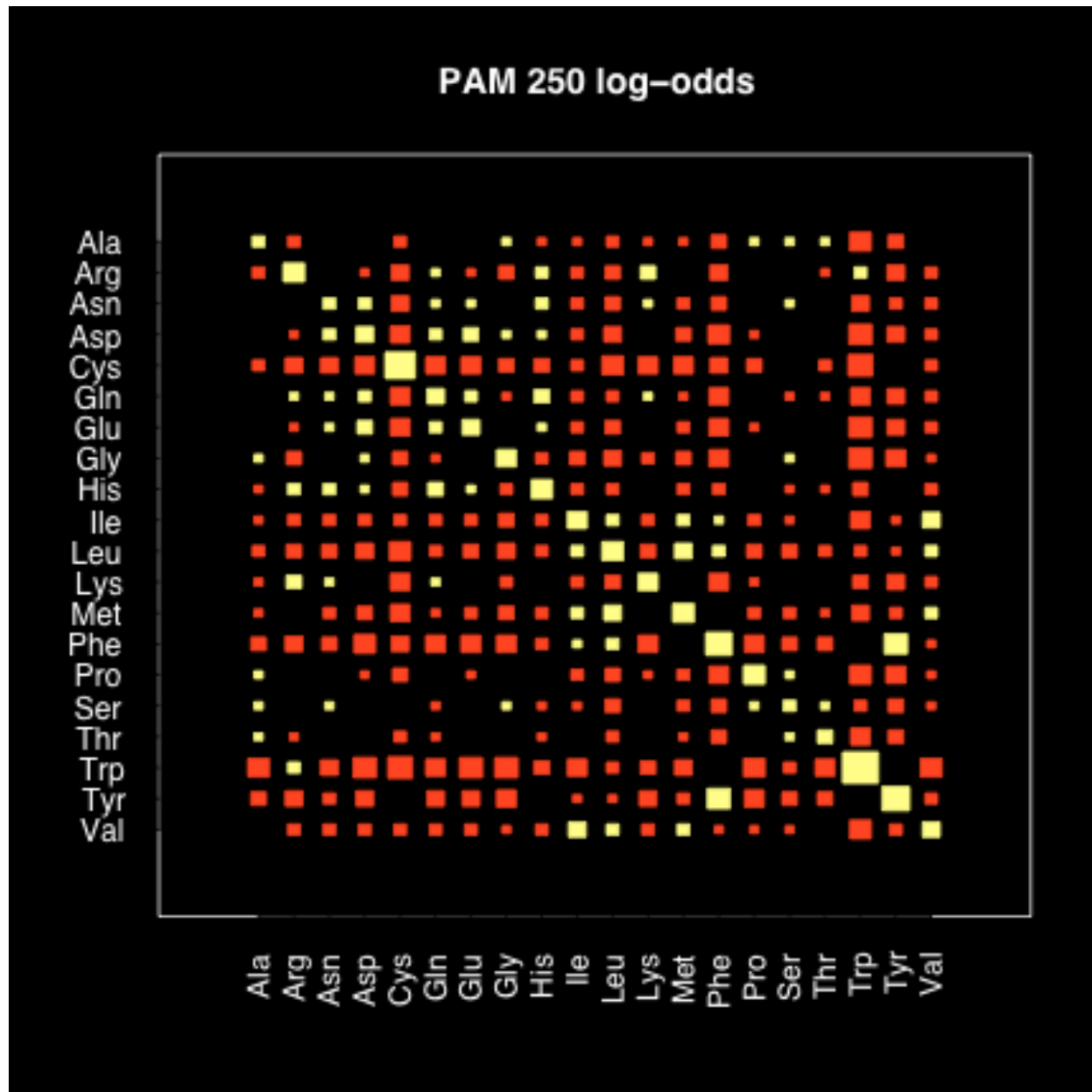
Extrapolation de la série de matrices PAM à partir de la PAM001

- Exemple: si l'on dispose de la matrice PAM001 (temps évolutif donnant ~1% de substitutions/position), on peut calculer la probabilité de substitution de l'Asn à la Thr en 2 unités temporelles (PAM002) en calculant la probabilité de chaque « trajet » de 2 substitutions.



$$\begin{aligned}
 P(\text{Asn} \rightarrow \text{Thr}) &= P(\text{Asn} \rightarrow \text{Ala} \rightarrow \text{Thr}) + P(\text{Asn} \rightarrow \text{Arg} \rightarrow \text{Thr}) + \dots + P(\text{Asn} \rightarrow \text{Val} \rightarrow \text{Thr}) \\
 &= (0.0009)(0.0001) + (0.0001)(0.0002) + \dots + (0.0001)(0.0009)
 \end{aligned}$$

Hinton diagram of the PAM250 matrix

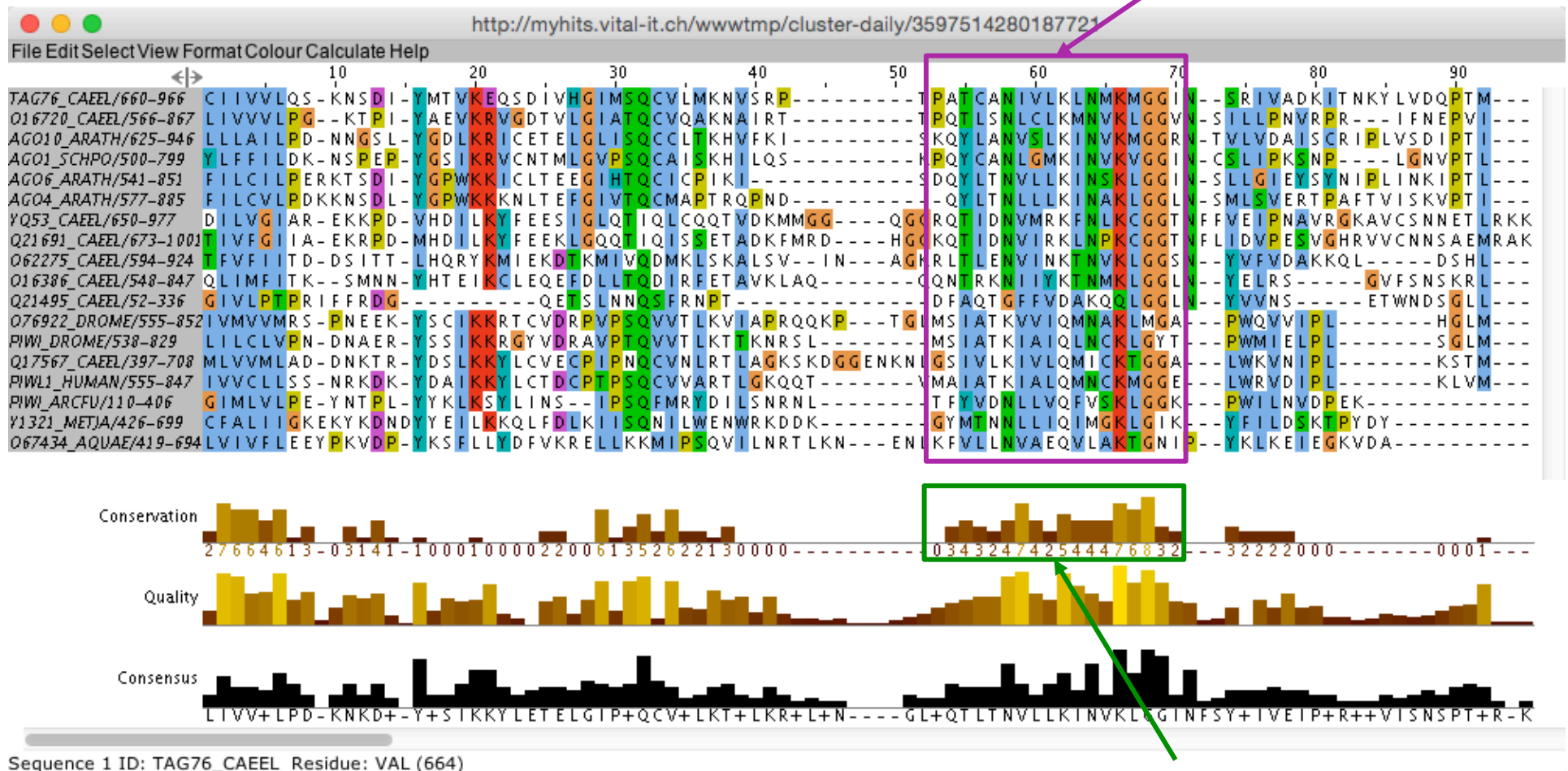


- Yellow boxes indicate positive values (accepted mutations)
- Red boxes indicate negative values (avoided mutations).
- The area of each box is proportional to the absolute value of the log-odds score.

La série BLOSUM
matrices de substitutions construites
à partir de blocs conservés

- Henikoff and Henikoff (1992) ont analysé les fréquences de substitutions dans des **blocs d'alignements multiples** générés à partir d'un grand nombre de familles de protéines (**blocks**).
- Ils en ont dérivé la série de matrices « **BLOSUM** », qui correspondent à des taux différents de **conservation évolutive** entre les séquences.

bloc d'alignement multiple



Taux de conservation

Sequence 1 ID: TAG76_CAEEL Residue: VAL (664)

BLOSUM scoring matrices

- Henikoff and Henikoff (1992) ont analysé les fréquences de substitutions dans des blocs d'alignements multiples générés à partir d'un grand nombre de familles de protéines (**blocks**)
- Ils en ont dérivé la série de matrices « **BLOSUM** », qui correspondent à des taux différents de divergence évolutive entre les séquences.
- Exemples
 - La matrice *BLOSUM62* a été calculée sur des blocs de $\geq 62\%$ d'identité
 - La matrice *BLOSUM80* a été calculée sur des blocs de $\geq 80\%$ d'identité
- Quand on utilise les matrices BLOSUM pour aligner des séquences, on devrait systématiquement **choisir la matrice la plus adéquate**, en fonction du pourcentage de similarité.
- Le problème est qu'avant de réaliser l'alignement, on connaît pas ce pourcentage. Comment résoudre cette circularité ?
 - On réalise un premier alignement avec une matrice « moyenne » (BLOSUM62).
 - On observe le % d'identité dans cet alignement.
 - On choisit alors la matrice dont l'indice est le plus proche de ce taux
 - On refait l'alignement avec la nouvelle matrice (sauf s'il s'agit de celle de départ).
- Exemples:
 - L'alignement présente 65.2% d'identité -> le premier alignement avec BLOSUM62 était correct.
 - L'alignement présente 28.4% d'identité -> on refait l'alignement avec BLOSUM30.
 - L'alignement présente 81.5% d'identité -> on refait l'alignement avec BLOSUM80.

***Utilisation des matrices de substitution
pour mesurer la qualité d'un alignement***

Matrices de substitutions nucléotidiques

- Pour les séquences nucléotidiques, on utilise généralement une pénalité identique pour toutes les substitutions.
- Cependant, on pourrait décider d'assigner un coût plus faible à certaines substitutions (par exemple A \leftrightarrow T) si l'on considère qu'elles ont plus de chance d'être observées dans des alignements (dans certains génomes, les résidus A et T sont deux fois plus fréquents que les C et G).
- Exemple: la matrice ci-jointe représente des scores définis de façon arbitraire
 - Identité +2
 - Substitution A-T -1
 - Autres substit. -2

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

Construction d'une matrice d'alignement à partir d'une matrice de substitutions

- Revenons à l'exemple précédent.
 - Haut: une matrice de scores (arbitraires) de substitution pour séquences d'ADN (un score pour chaque paire de nucléotides).
 - Bas: matrice d'alignement pour deux petites séquences d'ADN.
- Dans chaque cellule de la matrice d'alignement, on insère le score de la paire de résidus correspondants, extrait de la matrice de substitutions.

Matrice de substitutions

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

Matrice d'alignement

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

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

■ 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.