

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/>

Construction des matrices de substitutions

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.

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

Exemple d'alignement par paires

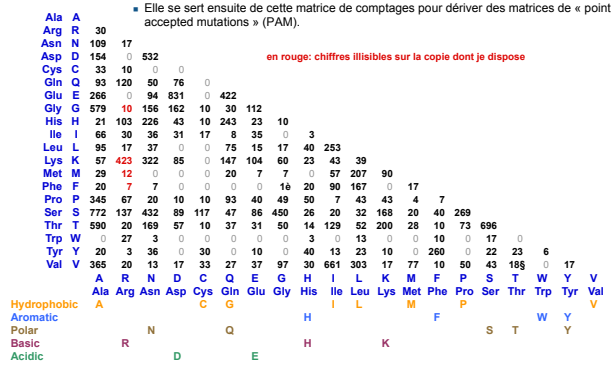
```
# Matrix: EHLGSUMS2
# Gap_penalty: 10.0
# Extend_penalty: 0.5
# Length: 482
# Identity: 133/482 (27.6%)
# Similarity: 205/482 (42.5%)
# Gaps: 85/482 (17.6%)
# Score: 353.5

metL 16 KFGQSLADVVCYTLRVAGIMAYTSQFDHDDHVVVAAGSTTQQLINMK-L-S 64
      ||||| | | | | | | | | | | | | | | | | | | | | | | | | | |
lysc  8 KFGQTSVADFDAMRRSADIVLSDANV-RLVVLSAAGILNLLVALLAELE 56
      . . . . . | . . . . . | . . . . . | . . . . . | . . . . . |
metL 65 QTDLSARVQVQTLRRYQCCLISGL----LFAEEADSLISAFVSOLESLA 110
      . . | . . . . . | . . . . . | . . . . . | . . . . . | . . . . . |
lysc  57 PQERF---EKLDATRHQIFAILERLKPFWKEEIERLEEN-IVYLAFAA 102
      . . . . . | . . . . . | . . . . . | . . . . . | . . . . . |
metL 111 ALLDGGINDAVVAFVQVQEVWASRLMSAVLNQQLPAMLALAREPFA-- 159
      ||| . | . . . | | | | | | | | | | | | | | | | | | | |
lysc  103 ALATS---PALTDELSYRDELMLTLPFVILIRERDVAQGFVQVREVMSTN 149
      . . . . . | . . . . . | . . . . . | | | | | | | | | | | | |
metL 160 ERAAGPVQVDSLEFLQGLLYVQPKRLVTV-GFISRNNAEYVLLGRN 208
      | . . . . | . . . . . | . . . . . | | | | | | | | | | | |
lysc  150 DRFGRAEFDIAALAEALQLPRINSLVITQGFISENKORTYTLGRG 199
      . . . . . | . . . . . | . . . . . | . . . . . | . . . . . |
metL 209 GSDYSATQIGALAGVSVTVTMSDVAGVYADRPKRVKDACLLPLRLDEAS 258
      ||| | | . . . . . | | | | | | | | | | | | | | | | | |
lysc  200 GSDYTAALAEALASAVLQVTVTKVYVTVVSVVSAACLEIPIAFVAAA 246
      . . . . . | . . . . . | . . . . . | . . . . . | . . . . . |
metL 259 ELKRLAAPLQARTLQVPSGSEIDQLRCVITFDQDSTRI-----E 299
      ||| . | | | | | | | | | | | | | | | | | | | | | |
lysc  250 EMATYQAVLAPATLLFAVRSDFPVPVSGSKDPRAAGTLVCKNTEHPPLF 299
      | . . . . . | | . . . . . | | | | | | | | | | | | |
metL 300 RVLASGTGAVITVSHDQVCLIEFPVPSAGDFKLAKKIDQLKRAQVPL 349
      | | | . . . . . | | . . . . . | | | | | | | | | | | |
lysc  300 KALALRRNQTLLTR-----SLNLSHRDF-LA--KVPILLAR----- 334
      . . . . . | . . . . . | . . . . . | . . . . . | . . . . . |
metL 350 AVGVNDHQQLGQFCYTSREVA-----DGL--KILDAGALPG 383
      | . . . . | | | | | | | | | | | | | | | | | |
lysc  335 ----NMS--VGLITTSREVSVALTDTGTSVNDVLLQSLMLLSALC 378
      . . . . . | . . . . . | . . . . . | . . . . . | . . . . . |
metL 384 ERLRNLQGLLVANVVGQVTR-----NPLKIRFVQQLQKQVPE 421
```

- La figure représente l'alignement de deux séquences peptidiques.
- Les barres verticales indiquent les identités.
- Les gaps sont marqués par des traits d'union.
- Les doubles points indiquent des substitutions qu'on retrouve souvent dans les alignements (« point accepted mutations »).
- Les simples points indiquent les substitutions rares et celles qui ne sont pas spécialement fréquentes.

Occurrences de substitutions dans 71 groupes de protéines alignées (Dayhoff, 1978)

- En 1978, Margret Dayhoff réalise des alignements de séquences protéiques (71 groupes de protéines), et compte le nombre de substitutions et d'identités entre chaque paire d'acides aminés.
- Elle obtient les comptages représentés dans la matrice ci-dessous.
- Elle se sert ensuite de cette matrice de comptages pour dériver des matrices de « point accepted mutations » (PAM).



Substitution matrices for proteins

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

	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.

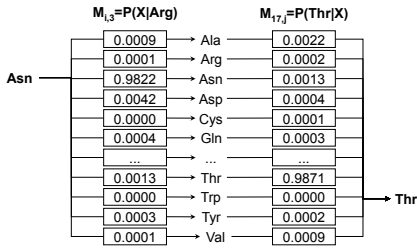
PAM scoring matrices

- 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.



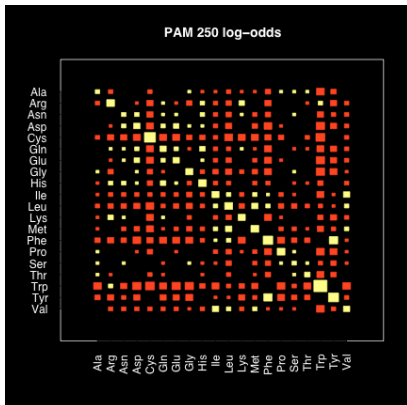
$$P(\text{Asn} \rightarrow \text{Thr}) = P(\text{Asn} \rightarrow \text{Ala} \rightarrow \text{Thr}) + P(\text{Asn} \rightarrow \text{Arg} \rightarrow \text{Thr}) + \dots + P(\text{Asn} \rightarrow \text{Val} \rightarrow \text{Thr})$$

$$= (0.0009)(0.0001) + (0.0001)(0.0002) + \dots + (0.0001)(0.0009)$$

PAM250 matrix

- La PAM250 est appropriée pour les alignements entre séquences très éloignées.
- Notes
 - La diagonale est constituée de scores positifs, qui reflètent la conservation.
 - les autres scores élevés correspondent souvent à des acides aminés partageant des propriétés physico-chimiques.

Hinton diagram of the PAM250 matrix



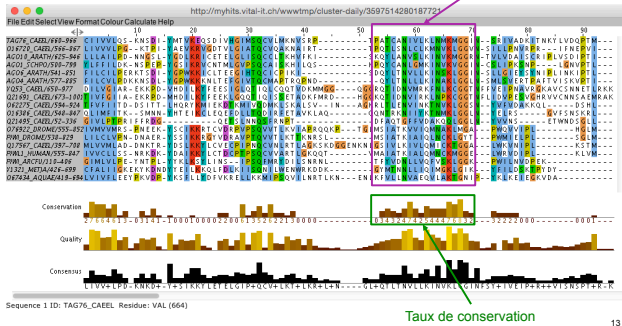
- Yellow boxes indicate positive values (accepted mutations)
- Red boxes indicate negative values (avoided mutations).
- The area of each box is proportional to the absolute value of the log-odds score.

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

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 **conservation évolutive** entre les séquences.

bloc d'alignement multiple



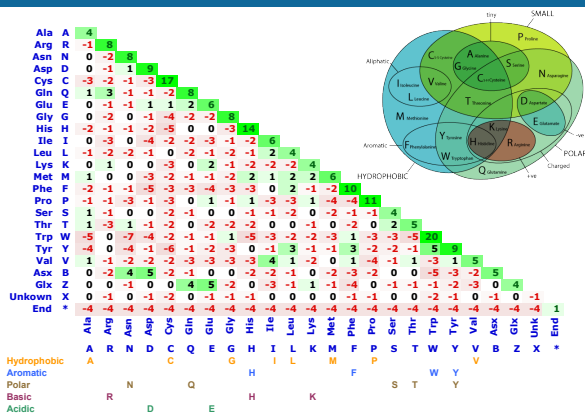
- Henikoff and Henikoff (1992) ont analysé les fréquences de substitutions dans des blocs d'alignements multiples générés à partir d'un grand nombre de familles de protéines (**blocks**).
- Ils en ont dérivé la série de matrices « **BLOSUM** », qui correspondent à des taux différents de divergence évolutive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

Reference: Henikoff, S. and Henikoff, J.G. (1992). Amino acid substitution matrices from protein blocks. PNAS 89:10915-10919.

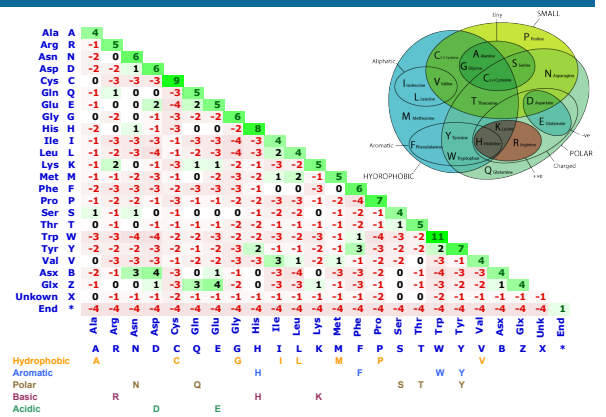
14

BLOSUM30



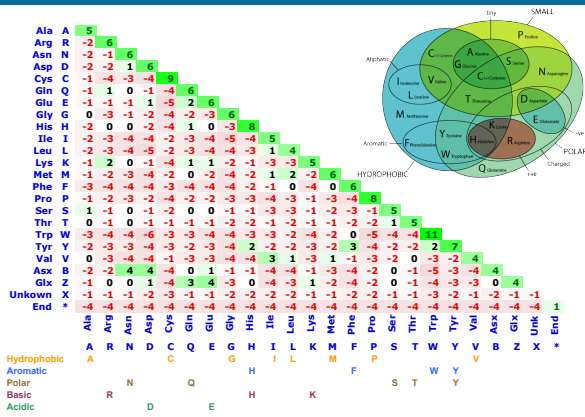
15

BLOSUM62



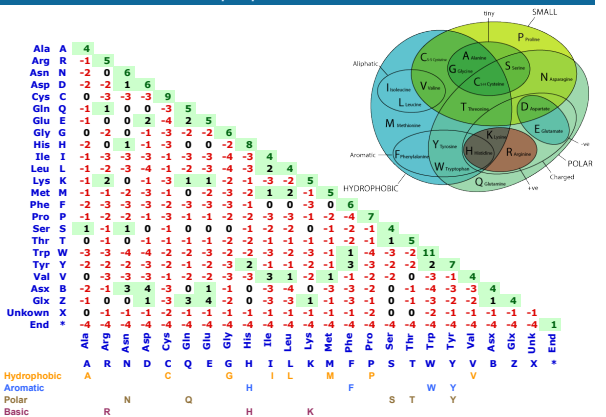
16

BLOSUM80



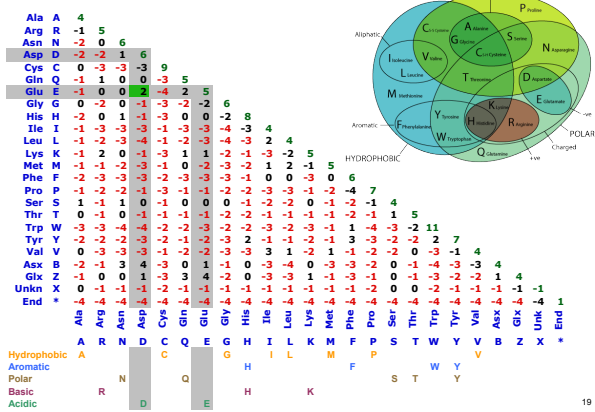
17

BLOSUM62 – Amino acid properties

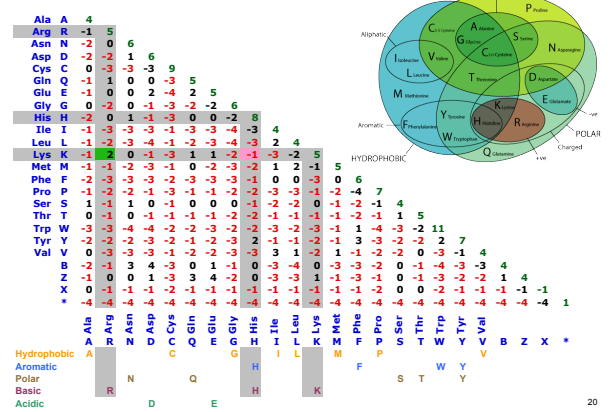


18

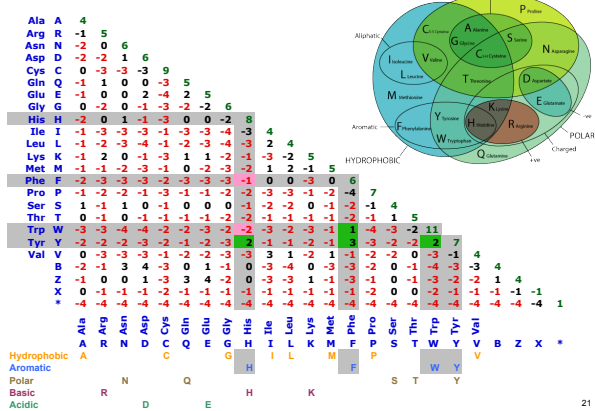
BLOSUM62 - substitutions between acidic residues



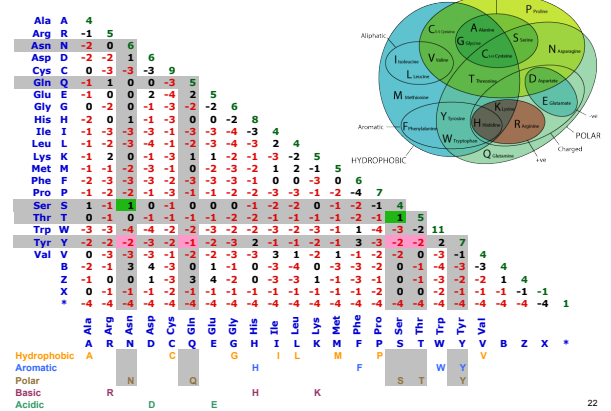
BLOSUM62 - substitutions between basic residues



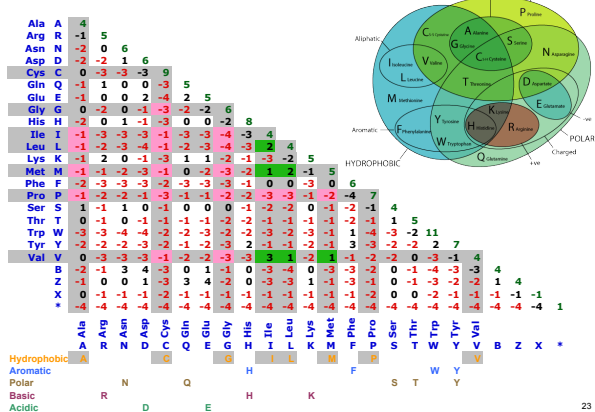
BLOSUM62 - substitutions between aromatic residues



BLOSUM62 - substitutions between polar residues



BLOSUM62 - substitutions between hydrophobic residues

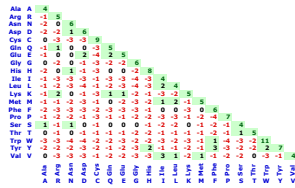


Utilisation des matrices de substitution pour mesurer la qualité d'un alignement

Matrices de substitutions

- Une **matrice de substitution** associe un score à chaque paire de résidus qu'on peut trouver dans un alignement.
 - Chaque ligne et chaque colonne représente l'un des résidus (4 nucléotides, 20 acides aminés).
 - La diagonale correspond aux identités.
 - Le triangle inférieur correspond à des substitutions.
 - Le triangle supérieur est symétrique au triangle inférieur, il n'est pas nécessaire d'indiquer les nombres.
- Les **scores négatifs** sont considérés comme des pénalités associées à certaines substitutions qu'on n'observe que rarement dans les alignements. Les algorithmes d'alignements tenteront donc d'éviter ces substitutions.
- Les **scores positifs** correspondent à des substitutions qu'on observe plus souvent que prévu, dans les alignements d'un grand nombre de séquences. Ceci suggère que ces substitutions particulières sont moins dommageable que d'autres, et on les qualifie donc de « **substitutions conservatives** » ou encore de « **mutations ponctuelles acceptées** » (**PAM**).
- Au sein d'un alignement, le terme **similarité** désigne les positions où se superposent des résidus ayant un score positif dans la matrice de substitution (identité ou substitution conservative).

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



25

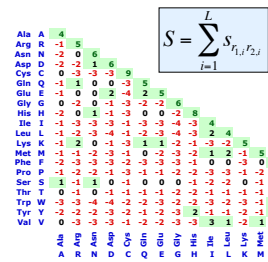
Matrices de substitutions nucléotidiques

- 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<->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

26

Utilisation d'une matrice de substitution pour calculer le score d'un alignement

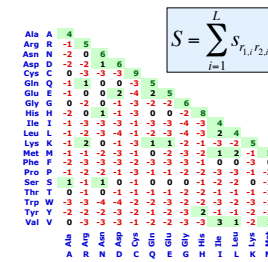


- Les matrices de substitution sont utilisées pour calculer le score d'un alignement.
- Ce score est la somme, pour toutes les positions de l'alignement (i de 1 à L), des scores des paires de résidus ($r_{1,i}$ et $r_{2,i}$).
- Les "gaps" sont traités par une règle spécifique reposant sur deux paramètres de pénalité:
 - Pénalité d'ouverture de gap (**go**)
 - Valeurs typiques: entre -10 et -15
 - Pénalité d'extension de gap (**ge**)
 - Valeurs typiques: entre -0.5 et -2

i	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21				
R	L	A	S	V	E	T	D	M	P	-	-	-	-	-	-	L	T	L	R	Q	H				
-	-	T	L	T	S	L	O	T	T	L	K	N	L	K	E	M	A	H	L	G	T	H			
S	-	1	4	0	0	4	1	2	5	-1	-2	-1	-10	-1	-1	-1	-1	-1	-1	-2	4	-2	-1	8	7

27

Utilisation d'une matrice de substitution pour calculer le score d'un alignement



- Les matrices de substitution sont utilisées pour calculer le score d'un alignement.
- Ce score est la somme, pour toutes les positions de l'alignement (i de 1 à L), des scores des paires de résidus ($r_{1,i}$ et $r_{2,i}$).
- Les "gaps" sont traités par une règle spécifique reposant sur deux paramètres de pénalité:
 - Pénalité d'ouverture de gap (**go**)
 - Valeurs typiques: entre -10 et -15
 - Pénalité d'extension de gap (**ge**)
 - Valeurs typiques: entre -0.5 et -2

i	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21				
R	L	A	S	V	E	T	D	M	P	-	-	-	-	-	-	L	T	L	R	Q	H				
-	-	T	L	T	S	L	O	T	T	L	K	N	L	K	E	M	A	H	L	G	T	H			
S	-	1	4	0	0	4	1	2	5	-1	-2	-1	-10	-1	-1	-1	-1	-1	-1	-2	4	-2	-1	8	7

28

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

29

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

30

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.