

Comparaisons locales et matrices de score

Équipe Bonsai
<http://www.lifl.fr/bonsai>

année 2013



Contenu du cours

- Matrices de scores
- Recherches locales : BLAST et FastA

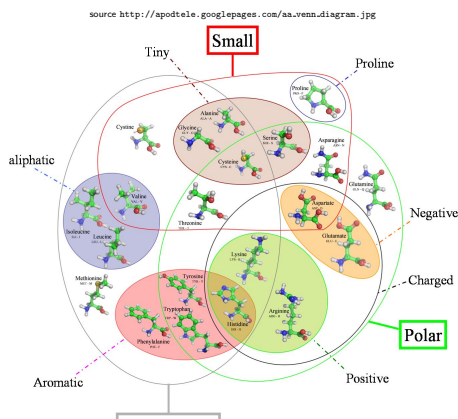
Exemples pour l'ADN

Identité (similarité)					BLAST (similarité)					Transition/Transversion				
	A	C	G	T	A	C	G	T	A	C	G	T		
A	1	0	0	0	1	-3	-3	-3	5	-4	-3	-4		
C	0	1	0	0	-3	1	-3	-3	-4	5	-4	-3		
G	0	0	1	0	-3	-3	1	-3	-3	-4	5	-4		
T	0	0	0	1	-3	-3	-3	1	-4	-3	-4	5		

Matrice BLOSUM-62 pour les protéines

```
# Matrix made by matlab from blosum62.ii
# * column uses minimum score
# BLOSUM Clustered Scoring Matrix in 1/2 Bit Units
# Blocks Database = /data/blocks.5.0/blocks.dat
# Cluster Percentage: = 62
# Entropy = 0.6979, Expected = -0.5209
A R N D C Q E G H I L K M F P S T V W Y V B Z X *
A -1 -2 -2 0 -1 -1 0 -2 -1 -1 -1 -1 -2 1 0 -3 -2 0 -2 1 0 -4
R 1 5 0 -2 -3 1 0 -2 0 -3 -2 -2 -1 -3 -2 -1 -1 -3 -2 -3 1 0 -1 -4
N -2 0 5 1 -3 0 0 0 1 -3 -3 0 -2 -3 -2 1 0 -4 -2 -3 3 0 -1 -4
D -2 -2 1 4 -3 0 2 -1 -1 -3 -4 -1 -3 -3 -1 0 -1 -4 -3 -3 4 1 -1 -4
C 0 -3 -3 -3 -3 -4 -3 -3 -1 -1 -3 -1 -2 -3 -1 -1 -2 -3 -1 -3 -2 -4
Q -1 1 0 0 -3 2 -2 2 0 -3 -2 -1 0 -3 -1 0 -1 -2 -1 -2 0 3 -1 -4
E -1 0 0 2 -4 -2 5 -2 0 -3 -3 1 -2 -3 -1 0 -1 -3 -2 -2 1 4 -1 -4
G 0 2 0 -1 -3 -2 2 -2 -4 -4 -2 -3 -3 -2 0 -2 -3 -3 -1 -2 -1 -4
H -2 0 1 -1 -3 0 0 -2 8 -3 -3 -1 -2 -1 -2 -1 -2 -2 2 -3 0 0 -1 -4
I -1 -3 -3 -3 -1 -3 -3 -4 -3 4 2 -3 1 0 -3 -2 -1 -3 -1 3 -3 -3 -1 -4
L -1 -2 -3 -4 -1 -2 -3 -4 -3 2 4 -2 2 0 -3 -2 -1 -4 -1 1 -4 -3 -1 -4
K -1 2 0 -1 -3 1 1 -2 -1 -3 -2 5 -1 -3 -1 0 -1 -3 -2 0 1 -1 -4
M -1 -1 -2 -3 -1 0 -2 -3 -2 1 2 -1 5 0 -2 -1 -1 -1 -1 -1 -3 -1 -1 -4
F -2 -3 -3 -3 -2 -3 -3 -3 -1 0 0 -3 0 6 -4 -2 -2 1 3 -1 -3 -3 -1 -4
P -1 -2 -2 -1 -3 -1 -1 -2 -2 -3 -3 -1 -2 -4 -1 -1 -4 -3 -2 -2 -1 -2 -4
S -1 -1 1 0 -1 0 0 0 -1 -2 -2 0 -1 -2 -1 1 1 -3 -2 -2 0 0 0 -4
T 0 -1 0 -1 -1 -1 -1 -2 -2 -1 -1 -1 -1 -2 -1 1 1 -2 -2 0 -1 -1 0 -4
W -3 -3 -4 -4 -2 -2 -3 -2 -2 -3 -2 -3 -1 1 -4 -3 -2 1 1 2 -3 -3 -2 -4
V -2 -2 -2 -3 -3 -1 -2 -3 2 -1 -1 -2 1 3 -3 -2 -2 2 7 -1 -3 -2 -1 -4
V 0 -3 -3 -3 -1 -2 -2 -3 -3 3 1 -2 1 -1 -2 -2 0 -3 -1 4 -3 -2 -1 -4
B -2 -1 3 4 -3 0 1 -1 0 -3 -4 0 -3 -3 -2 0 -1 -4 -3 -3 4 1 -1 -4
Z -1 0 0 1 -3 3 4 -2 0 -3 -3 1 -1 -3 -1 0 -1 -3 -2 2 1 4 -1 -4
X 0 -1 -1 -1 -2 -1 -1 -1 -1 -1 -1 -1 -1 -1 -2 0 0 -2 -1 -1 -1 -1 -4
* -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4 -4
```

Matrice BLOSUM-62 pour les protéines



De l'importance des matrices de scores

- impliquées dans toutes les analyses par comparaison de séquences
- résultats **fortement** dépendants de la matrice
- représentent implicitement une théorie de l'évolution (matrices protéiques)
- la compréhension d'une matrice ⇒ un bon choix

Similarité vs. distance

- un élément de la matrice représente :
 - le coût du remplacement d'une base par une autre (*distance*)
 - la mesure de la *similarité* du remplacement
- association entre
 - distance → phylogénie
 - similarité → recherche dans des bases de données
- même principe de recherche :
 - **maximiser un score** ≡ **minimiser une distance**
 - matrice de distance et de similarité peuvent être déduites l'une de l'autre

Comment obtenir une telle matrice ?

pour l'ADN
 ⇒
 souvent données de manière *ad hoc*

Comment obtenir une telle matrice ?

pour les protéines

- matrices *log odds ratio*

$$S_{ij} = \log \frac{q_{ij}}{p_i p_j}$$
- exprime le ratio entre :
 - la probabilité que deux résidus *i* et *j* soient alignés par **descendance**
 - et la probabilité que ceux-ci soient alignés par **chance**
- explication :
 - q_{ij} = la fréquence que l'alignement de *i* et *j* soit observé dans des séquences homologues.
 - p_i = la fréquence d'occurrence de *i*
 - un score est > 0 si la proba d'un match significatif est > à la proba d'un match aléatoire

⇒ matrices PAM et BLOSUM

Matrice BLOSUM

Blocks **S**ubstitutions **M**atrices
 Henikoff & Henikoff, 1992

- fréquence de changements entre deux acides aminés avec conservation de structure
échantillon : BLOCKS (alignements multiples sans gaps)
- BLOSUM-*N* : seuil de similarité, *N* = % de similarité
- convient bien pour la recherche de similarités locales
- la plus courante : BLOSUM-62 (matrice par défaut de BLAST)

Matrice BLOSUM - Matériel

Matrice BLOSUM - Schéma de Construction

- 1 éliminer les alignements d'identité < *n*% ⇒ **BLOSUM-*n***.
- 2 compter le nombre f_{ij} de paires d'aa *i* et *j*.
- 3 calcul de la fréquence q_{ij} des paires d'aa *i* et *j* :

$$q_{ij} = \frac{f_{ij}}{\# \text{ total de paires}}$$

- 4 calcul de la fréquence *marginale* p_i des aa *i* :

$$p_i = \sum_j q_{ij}$$

(plus exactement $p_i = q_{ii} + \sum_{j \neq i} q_{ij}/2$)

- 5 calcul de la matrice log odds S_{ij} :

$$S_{ij} = c \times \log \frac{q_{ij}}{p_i p_j}$$

(plus exactement $S_{i \neq j} = c \times \log \frac{q_{ij}}{2p_i p_j}$ et $S_{ii} = c \times \log \frac{q_{ii}}{p_i p_i}$)

Matrice PAM

Accepted Point Mutation / Percent Accepted Mutation
Dayhoff, 1979

- fréquence de changements entre acides aminés
Reconstitution de l'évolution avec la construction d'arbres phylogénétiques pour 71 familles de protéines
- convient bien pour les séquences avec un ancêtre commun
- possibilité de choix d'une matrice en fonction de l'évolution supposée
PAM-N : N mutations acceptées par 100 acides aminés
- si la distance mutationnelle n'est pas connue, faire plusieurs essais
PAM 40, PAM 120, PAM 250, par exemple.

Matrice PAM - Construction

- alignements globaux de familles de protéines (identité > 85%)
- reconstruction d'une phylogénie et des ancêtres (71 familles)
- compter le nombre de fois A_{ij} où un aa i est remplacé par un aa j dans toutes les comparaisons 2 à 2
- estimation de la mutabilité m_j d'un aa j
- calcul de la matrice de probabilité de mutations

$$M_{ij} = \lambda \frac{m_j A_{ij}}{\sum_i A_{ij}} \text{ et } M_{jj} = 1 - \lambda m_j$$

m_j : probabilité pour j de muter

$\frac{A_{ij}}{\sum_i A_{ij}}$: probabilité conditionnelle pour j , s'il mute, de muter en i

- calcul de la matrice log odds

$$S_{ij} = \log \frac{M_{ij}}{p_i}$$

$S_{ij} = \log \frac{p_i \times M_{ij}}{p_i p_j}$: comme p_j est la proba de j dans une séq., et M_{ij} la proba, pour un occ. j , de muter en i durant un laps de temps donné, le

Matrice PAM - Construction (1/6)

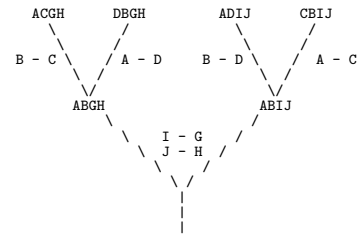
alignements de 71 groupes de l'Atlas of protein sequence

```

KAPPA
1 HUMAN EU
2 MOUSE WGPC 21
3 CAT 2211
4 RABBIT 4135
5 RABBIT
6 HUMAN SH
7 PIG
8 2 MOUSE WGPC 104E
9 2 MOUSE WGPC 315
BETA-2 MICROGLOBULIN
10 HUMAN
HEAVY CHAIN FIRST
11 GAMMA-1 EU
12 BETA-2 MICROGLOBULIN
13 HUMAN BUR
14 HUMAN BUR
HEAVY CHAIN EXTRA
15 BETA-2 MICROGLOBULIN
16 BETA-2 MICROGLOBULIN
HEAVY CHAIN MIDDLE
17 GAMMA-1 EU
18 BETA-2 MICROGLOBULIN
19 HUMAN BUR
20 HUMAN BUR
HEAVY CHAIN LAST
21 GAMMA-1 EU
22 BETA-2 MICROGLOBULIN
23 ALPHA-1 BUR
24 HUMAN BUR
CONSERVED
P V P L C L V G F P V V W
    
```

Matrice PAM - Construction (2/6)

phylogénie et reconstruction des ancêtres pour chacun des 71 groupes



Matrice PAM - Construction (3/6)

nombre d'accepted point mutations

	ala	arg	asn	asp	cys	gln	glu	gly	his	ile	leu	lys	met	phe	pro	ser	thr	trp	tyr	val
ala	30																			
arg	109	17																		
asn	154	0	532																	
asp	33	10	0	0																
cys	93	120	50	76	0															
gln	266	0	94	331	0	422														
glu	579	10	156	162	10	30	112													
gly	21	103	226	43	10	243	23	10												
his	66	30	36	13	17	8	35	0	3											
ile	95	17	37	0	0	75	15	17	40	253										
leu	57	477	322	85	0	147	104	60	23	43	39									
lys	29	17	0	0	0	20	7	7	0	57	207	90								
met	20	7	0	0	0	17	20	50	157	0	17									
phe	345	67	27	10	10	93	40	49	50	7	43	43	4	7						
pro	772	137	432	98	117	47	86	450	26	20	32	168	20	40	269					
ser	590	20	169	57	10	37	31	50	14	129	52	200	28	10	73	696				
thr	0	27	3	0	0	0	0	3	0	13	0	0	10	0	17	0				
trp	0	3	36	0	30	0	10	0	40	13	23	10	0	260	0	22	23	6		
tyr	365	20	13	17	33	27	37	97	30	661	303	17	77	10	50	43	186	0	17	
val																				

Matrice PAM - Construction (4/6)

Exemple de calcul :

Alignement	A D A	A D B
Acides aminés	A B D	
Chgt. observés	1 1 0	
Fréquence d'occ.	3 1 2	
Mutabilité relative	.33 1 0	

Mutabilités relatives de tous les aa :

Ser	149	Ala	100	Gln	98
Met	122			Asp	90
Asn	111			Thr	90
Ile	110			Gap	84
Glu	102			Val	80
				Lys	57
				Pro	56
				His	50
				Gly	48
				Phe	45
				Arg	44
				Leu	38
				Tyr	34
				Cys	27
				Trp	22

FASTA

Pearson et Lipman, 1988

- alignement **global** avec gaps
- traite les séquences de la banque les unes après les autres
- fonctionnement :
 - 1 trouve tous les mots identiques de longueur $\geq l$ communs à q et t_i
 - 2 sélectionne ceux de score suffisamment élevé (score PAM par exemple)
 - 3 sélectionne une diagonale d (du dotplot) contenant le maximum de mots identiques de longueur $\geq l$
 - 4 procède à un alignement global "classique" dans une bande de largeur $2k$ autour de la diagonale d
- deux paramètres : k et l , l généralement de longueur 6 pour l'ADN et 2 pour les protéines

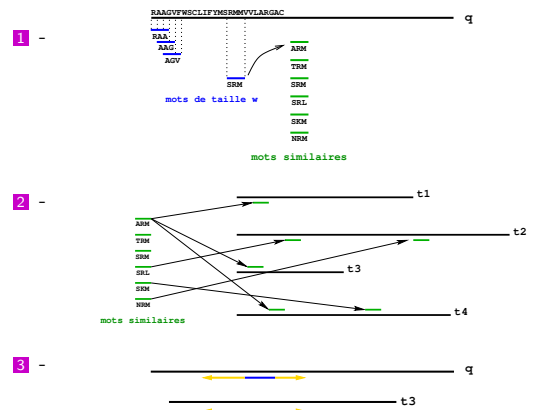
Blast

- naît en 1990 : trouve des matchs significatifs sans gaps
- évolution vers une version 2, avec gaps
 - NCBI-Blast
 - WU-Blast : très similaire à NCBI-Blast (mix entre Blast1 et FASTA pour la dernière étape)
- évolution vers des versions avec raffinement des résultats

Blast 1

- Trouve les mots **similaires** de taille fixe entre q et t_i .
(taille par défaut : ADN $\rightarrow 11$, Protéines $\rightarrow 3$)
- Ne considère que les couples de mots ...
- (Protéines) **similaires** \rightarrow score des mots alignés \geq seuil T
(à l'origine $T = 13$, actuellement $T = 11$ sur BLOSUM-62)
- (ADN) **identiques** \rightarrow pas de seuil T donc **moins sensible**.
- Chaque couple de mots entre q et un t_i forme un **hit**
- Chaque hit est étendu à gauche et à droite : l'extension est stoppée lorsque le score du hit décroît de plus de X (**X-drop**)

Blast 1 - Schématiquement



Blast 1

- un **hit** est un mot "commun" de taille fixée w (et de score supérieur à un seuil T dans le cas de BLAST-P) sur les deux séquences q et t_i
- chaque hit **étendu** (X-drop) forme un **LMSP** : *Locally Maximal scoring Segment Pair*
- ne conserve que les LMSP de score supérieur à un score seuil donné, les **HSP** : *High scoring Segment Pair*
- significativité évaluée (pour chaque t_i) sur le meilleur HSP trouvé nommé **MSP** : *Maximur scoring Segment Pair*

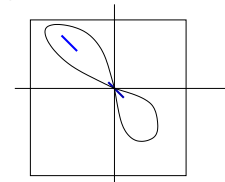
NCBI - Blast 2 (Gapped-blast)

idée 1 : se baser sur 2 hits distants au maximum de A sur la même diagonale (BLASTP)



\rightarrow baisser le seuil de score de chaque hit $T = 13 \rightarrow 11$ pour conserver une bonne sensibilité

idée 2 : étendre les hits comme dans Blast 1 (X-drop) mais en **autorisant les gaps**



MegaBLAST

pour l'ADN

- idée : un Blast plus rapide lorsqu'on recherche une grande similarité
- mise en œuvre : utiliser des mots de taille plus grande (28 contre 11)
- à réserver à des requêtes du style : trouver la séquence dans la banque
- évolution : **Discontiguous MegaBLAST**
 - principe : utiliser une *graine espacée* plutôt qu'un *mot exact* (graine contiguë)
 - exemple : graine espacée 100101100101100101101 plutôt que graine contiguë 11111111111
 - peut se révéler meilleur que BLAST (en particulier avec *graines espacées multiples*).

Définition des graines contiguës vs espacées

graine contiguë : 111111

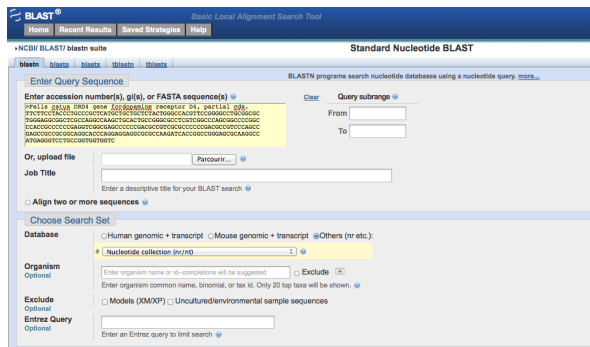
graine espacée : 11101011

```
111111
ATCAGTGC AAAATGCCGAAGA
|||||:|||||.|||||
ATCAGCGCAAATGCTCAAGA
```

```
11101011
ATCAGTGC AAAATGCCGAAGA
|||||:|:|||||.|||||
ATCAGCGCAAATGCTCAAGA
```

- Les graines espacées peuvent être *bien choisies* pour mieux détecter les alignements (*Keich Li Ma Tromp DAM 2004*)
- Il est possible d'utiliser plusieurs graines espacées de formes différentes pour améliorer la sensibilité de la recherche

Exemple : Basic BLAST → nucleotide blast → nr



Exemple : Résultats

```
Query= Felis catus DRD4 gene fordopamine receptor D4
(276 letters)
Database: All GenBank+EMBL+DDBJ+PDB sequences
1,174,453 sequences; 5,001,591,585 total letters

Sequences producing significant alignments:

Score E
(bits) Value
gi|AB069665 Felis catus DRD4 gene f... 210 5e-52
g|AB069662 Nyctereutes procyonoides... 157 7e-36
gi|AB069661 Canis lupus DRD4 gene f... 157 7e-36
gi|AB069666 Bos taurus DRD4 gene fo... 143 1e-31
gi|291947 Homo sapiens Dopamine D4 recep... 135 2e-29
```

Exemple : Résultats

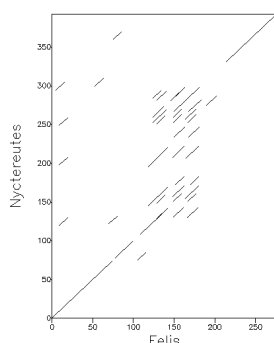
```
>gi|18143632|dbj|AB069662.1|AB069662 Nyctereutes procyonoides
DRD4 gene fordopamine receptor D4. Length = 393
Score = 157 bits (79), Expect = 7e-36
Identities = 94/99 (94%)
Strand = Plus / Plus
Query 1 ttcttctaccctgccgctcatgctgctgctctactggccagttcc 48
|||||:|||||:|||||:|||||:|||||:|||||:|||||:|||||:
Sbjct 1 ttcttctaccctgccgctcatgctgctgctctactggccagttcc 48
Query 49 ggggcctgaggcctggagggcggctcgccaggccaagctgcactcgg 99
|||||:|||||:|||||:|||||:|||||:|||||:|||||:|||||:
Sbjct 49 ggggcctgaggcctggagggcggcgtcgccaggccaagctgcactcgg 99
Score = 107 bits (54), Expect = 5e-21
Identities = 60/62 (96%)
Strand = Plus / Plus
Query 215 gtaggcggccaagatcaccggcgggagcgcaagcctgaggtcct 252
|||||:|||||:|||||:|||||:|||||:|||||:|||||:|||||:
Sbjct 332 gtaggcggccaagatcaccggcgggagcgcaagcctgaggtcct 379
Query 253 tgccggtggtgct 276
|||||:|||||:|||||:|||||:|||||:|||||:|||||:|||||:
Sbjct 380 tgccggtggtgct 393
```

BLAST : ... vs Alignement Réel

```
Felis Catus/ Nyctereute
1 ttcttctaccctgccgctcatgctgctgctctactggccagc 145 ggcgagc.....
1 ttcttctaccctgccgctcatgctgctgctctactggccagc 181 ggc agccggagcgaaccccgccgcccggccgagcgac
46 ttccgggctcggcctggagggcggctcgccaggccaagctg 152 .....
46 ttccgggctcggcctggagggcggcgtcgccaggccaagctg 225 cccgatgacaccccgagcgaacccctgcccccggcccccgc
91 cactgcgggcctcgtcgccagcggcccgccgcccagccgccc 153 cccgagcgcctgcccggcccccggccgctcccgagcggagccgc
91 cagggcgagacacggcagaccagcggcccgccgcccagccccc 270 cccgagcgcctcggcccccggccgcccggccgacctggagccag
136 ccga.ggt.....c 198 gggcagcaccacaggaggggcggcggccaagatcacggcggga
136 ccgaggttaccggcccccggcccccggcagcggccggcagc 315 gggcagcaccacaggaggggcggcggccaagatcacggcggga
243 gggcagggcctgagggtcctgcccgtggtggtc
360 gggcagggcctgagggtcctgcccgtggtggtc
```

BLAST : ... vs Alignement Réel

Felis Catus/ Nyctereute



BLAST : Significativité des MSPs

- deux séquences peuvent toujours être alignées
- il existe toujours un (au moins) alignement de meilleur score S entre deux séquences (un MSP)

question : ce score est-il suffisamment élevé pour prouver une homologie ?

problème : peut-on trouver un MSP de meilleur score dans deux séquences aléatoires ?

BLAST : Mesures de significativité

■ la p-valeur (p-value)

mesure la *Probabilité* que 2 séquences aléatoires de même longueur et de même composition possèdent un MSP de score $\geq S$

■ la E-valeur (E-value)

mesure l'*Esperance* E du nombre n de MSPs de score $\geq S$ dans 2 séquences aléatoires de même longueur et de même composition

$$E = \sum_n p(n) \times n$$

Calcul de la E-value

- soient deux séquences a et b aléatoires suivant une distribution de probabilité connue
- on suppose que les MSPs sont données par les diagonales du dotplot
- plutôt que de décrire un alignement par des paires de lettres tirées aléatoirement, on peut le décrire par une suite de scores tirés aléatoirement
- on veut calculer l'esperance du nombre de MSPs de score $\geq S$

Calcul de la E-value

- Selon *Karlin et Altschul 1991* :

$$\mathbf{E\text{-value}} = Kmn e^{-\lambda S} \quad \mathbf{p\text{-value}} = 1 - e^{-\mathbf{E\text{-value}}}$$

avec m la taille de la séquence requête, n la taille de la banque de données, S le score du hit (K et λ dépendent de la matrice de score, K peut être ajusté en fonction du coût des gaps)

Calcul du bit-score

- si S est le score d'un hit
- le bit-score (score normalisé) est :

$$S' = \frac{\lambda S - \ln K}{\ln 2}$$

- l'expression de la E-value devient :

$$E(S) = mn2^{-S'}$$

Variation de la E-value

- si la taille de la séquence query augmente : la E-value
- si la taille de la banque est divisée par deux : la E-value
- si le score augmente : la E-value
- quel bit-score pour obtenir une E-value de 0.05 pour une séquence de longueur 250 et une bd de longueur 50000000 ?
- si on passe la E-value à 0.01, quel sera le bit-score ?

Variations de la E-value

```
Mus musculus chromosome 5, clone RP23-30119, complete sequence
Length = 212246

Score = 32.2 bits (16), Expect = 2.1
Identities = 19/20 (95%)
Strand = Plus / Plus

Query: 2      ttcattatgaagcagcagga 21
             |||          |||
Sbjct: 136843 ttcattatgatgcacgagga 136862

Mus musculus BAC clone RP23-13L19 from chromosome 9, complete sequence
Length = 224108

Score = 30.2 bits (15), Expect = 8.1
Identities = 15/15 (100%)
Strand = Plus / Plus

Query: 6      ttatgaagcagcagg 20
             |||          |||
Sbjct: 93798  ttatgaagcagcagg 93812

Lambda      K      H
1.37        0.711  1.31
Number of Hits to DB: 99,084,306
Number of Sequences: 2130505
Number of extensions: 2852
Number of successful extensions: 19
Number of sequences better than 10.0: 0
Number of HSP's better than 10.0 without gapping: 0
Number of HSP's successfully gapped in prelim test: 16
length of query: 21
length of database: 10,249,863,584
```

Variations de la E-value

```
Mus musculus chromosome 5, clone RP23-30119, complete sequence
Length = 212246

Score = 32.2 bits (16), Expect = 2.1
Identities = 19/20 (95%)
Strand = Plus / Plus

Query: 2      ttcattatgaagcagcagga 21
             |||          |||
Sbjct: 136843 ttcattatgatgcacgagga 136862

-----

Mus musculus, clone RP23-277C24, complete sequence
Length = 199946

Score = 32.2 bits (16), Expect = 2.1
Identities = 16/16 (100%)
Strand = Plus / Minus

Query: 1      attcattatgaagcac 16
             |||          |||
Sbjct: 69080  attcattatgaagcac 69065
```

Variations de la E-value

```
Query length : 21
Mus musculus, clone RP23-277C24, complete sequence
Length = 199946

Score = 32.2 bits (16), Expect = 7.6
Identities = 16/16 (100%)
Strand = Plus / Minus

Query: 1      attcattatgaagcac 16
             |||          |||
Sbjct: 69080  attcattatgaagcac 69065

-----

Query length : 20
Mus musculus, clone RP23-277C24, complete sequence
Length = 199946

Score = 32.2 bits (16), Expect = 5.1
Identities = 16/16 (100%)
Strand = Plus / Minus

Query: 1      attcattatgaagcac 16
             |||          |||
Sbjct: 69080  attcattatgaagcac 69065
```

Les différents programmes BLAST

Query \ Database	nucléique	protéique	nucléique traduit
nucléique	blastn	x	x
protéique	x	blastp	tblastn
nucléique traduit	x	blastx	tblastx

Le bon programme pour la bonne requête

extrait de "BLAST Program Selection Guide"

- MEGABLAST is the tool of choice to identify a nucleotide sequence
- Discontiguous MEGABLAST is better at finding nucleotide sequences similar, but not identical, to your nucleotide query
- les pages "Search for short nearly exact matches"
 - nucleotide : useful for primer or short nucleotide searches
 - proteins : optimized to find matches to a short peptide
 - principales différences :
 - taille de mots plus petite
 - suppression des filtres
 - relâchement de la E-value
 - matrice de score PAM30 (au lieu de BLOSUM62) pour les protéines

Evolutions de Blast : PSI-Blast

PSI-Blast is designed for more sensitive protein-protein similarity searches

Position Specific Iterated BLAST

- 1 recherche initiale avec BLASTp
- 2 construction d'un alignement multiple, puis d'un profil
 - à partir d'un alignement multiple des meilleurs hits
 - construit une matrice position-spécifique :
 - chaque colonne représente un AA
 - chaque ligne une position dans l'alignement

3 nouvelle recherche avec le profil et modification réitérè le processus un certain nombre de fois ou jusqu'à convergence

Profil - exemple

```
Alignement multiple: # Pure Frequency Matrix
# Columns are amino acid counts A->Z
# Rows are alignment positions 1->n
Simple
T-VAAPSVFIFPPSDEQ Name mymatrix
A-DAAPTWSIFPPSSEQ Length 17
A-NAAPTWSIFPPSTZZ Maximum score 60
D-PVAPTVLIFPPAADQ Thresh 75
DFPIAPTLLFPPSADQ Consensus APPAAPTVLIFPPSADQ
200200000000000000000000000000000000
0000000000000000000001000000000000
0001000000000102000001000000
30000001000000000000000001000000
5000000000000000000000000000000000
0000000000000000050000000000000000
00000000000000000000014000000000
0000000000000000000000050000000000
0000100000200000020000000000000000
0000000040010000000000000000000000
0000050000000000000000000000000000
0000000000000000050000000000000000
0000000000000000050000000000000000
1000000000000000000000400000000000
2001000000000000000001100000000000
000220000000000000000000000000010
00000000000000000040000000000000
```

Evolution de Blast : PHI-Blast

PHI-BLAST can do a restricted protein pattern search

Pattern Hit Initiated BLAST

- pour les séquences protéiques
- entrée : une séquence et un motif (expression régulière à la Prosite)
- restriction de la banque aux séquences pour lesquelles le motif est retrouvé
- puis application de BLAST
- couplage possible avec PSI-Blast

Sensibilité et spécificité

On peut classer les résultats d'une méthode en 4 catégories :

VP les vrais positifs (*classé positif et bien positif*)

FP les faux positifs (*classé positif mais réellement négatif*)

ex. alignements parasites, prédiction pour une fonction que la séquence n'a pas en réalité...

FN les faux négatifs (*classé négatif mais réellement positif*)

ex. alignements perdus, prédiction négative pour une fonction que la séquence a en réalité...

VN les vrais négatifs (*classé négatif et bien négatif*)

$$\text{sensibilité} = \frac{VP}{VP+FN} \quad \text{spécificité} = \frac{VN}{VN+FP}$$

- sensibilité : capacité de la méthode à ne pas "louper" de Positifs
- spécificité : capacité de la méthode à ne pas "ramener" de Négatifs