

ALIGNEMENT DE DEUX SÉQUENCES

UNE PETITE FAUTE

Remarque : La visite de suivi effectuée par l'auditeur/l'inspecteur responsable à la suite d'une non-conformité majeure doit avoir lieu le plus rapidement possible après la « date d'exécution de l'action corrective » indiquée à la partie B de la DAC (voir la section 6.3.3 du présent document).

Lorsque l'auditeur/l'inspecteur responsable constate que l'action corrective a été exécutée et qu'elle est efficace, la DAC peut être classée.

(Source: Agence canadienne d'inspection des aliments. Manuel de mise oeuvre du Programme d'amélioration de la salubrité des aliments.)

«document» : document, dorment, doucement, ...

DISTANCE ENTRE DEUX MOTS

Distance d'édition : nombre d'opérations pour transformer un mot à un autre.

Opérations : sur caractères

- insertion (**I** *insert*)
- suppression (**D** *delete*)
- substitution (**S** *substitute*)
- identité (**M** *match*)

CALCUL DE DISTANCE

But : calculer le nombre *minimal* d'opérations pour la transformation.

Minimal : autant de **M** (identité) que possible.

Exemples : docment-document, docment-dorment, docment-doucement, docment-moments

On peut avoir plusieurs suites de la même taille minimale : abbc-axc

CALCUL DE DISTANCE 2

Idée : calculer les distances entre les **préfixes** successivement

Nos notions formelles :

- **alphabet** fini Σ (p.e. $\Sigma = \{A, \dots, z\}$ ou $\Sigma = \{T, G, A, C\}$)
- **mot** ou **séquence** : une suite de caractères de Σ

Soit S une séquence.

- **taille** de S , denotée par $|S| =$ nombre de caractères
- caractère en position i : $S[i]$
- **sous-mot** $S[i..j]$: le mot formé par les caractères $S[i], S[i + 1], \dots, S[j]$
- **préfixe** : sous-mot de forme $S[1..i]$.

CALCUL DE DISTANCE 3

Def. Distance d'édition entre deux séquences S_1 et S_2 : nombre minimal de **I**, **D**, **S** dans une suite d'opérations qui transforme S_1 en S_2 .

Thm. La distance d'édition est une fonction symétrique.

Preuve. **I** \Leftrightarrow **D**.

Thm. La distance d'édition satisfait l'inégalité du triangle



Un fait intéressant :) la distance d'édition est parfois appelée *distance de Levenshtein*

CALCUL DE DISTANCE 4

Def. Soit S et T deux séquences. On définit $D(i, j)$ par

$$D(i, j) = \begin{cases} \text{distance entre } S[1..i] \text{ et } T[1..j] & \text{si } i > 0 \text{ et } j > 0, \\ i & \text{si } j = 0 \\ j & \text{si } i = 0. \end{cases}$$

Donc $D(i, j)$ est la distance entre les deux préfixes de tailles i et j .

Thm. Si $i, j > 0$, on a

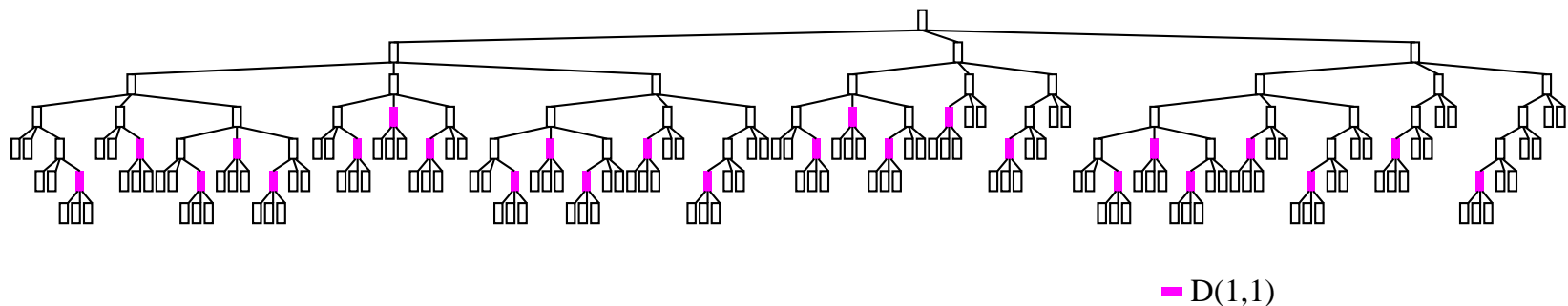
$$D(i, j) = \min \left\{ \begin{array}{l} D(i - 1, j) + 1, \\ D(i, j - 1) + 1, \\ D(i - 1, j - 1) + \{S[i] \neq T[j]\} \end{array} \right\}.$$

CALCUL DE DISTANCE 5

- Preuve.** 1. La dernière opération pour achever $D(i, j)$ doit être **I**, **D**, ou **S/M**.
2. Tous les trois sont possibles.

Donc on peut calculer $D(|S|, |T|)$ par récursion... pas une bonne idée. (La même valeur sera calculée plusieurs fois.)

Distance: opra - uma



Quand même on n'a que $(1 + |S|) \times (1 + |T|)$ appels récursifs possibles.

Au lieu d'explorer l'arbre de récursions en descendant, il est mieux de le faire de manière ascendante.

TABLEAU DE CALCULS

Les cases contiennent les $D(i, j)$.

Parcours : ligne par ligne.

Exemple : AAAC \rightarrow AGC

Temps de calcul : $O(mn)$ (on remplit chaque case du tableau en un temps constant : trois comparaisons et deux ou trois additions)

Comment trouver une suite d'opérations qui correspond à $D(i, j)$? Enregistrer dans chaque case la direction du min de la recurrence.

MAIS POURQUOI ÇA NOUS INTÉRESSE ?

L'alignement de séquences est une méthode utilisée très souvent en biologie moléculaire d'aujourd'hui.

Idée : similarité de séquences \Rightarrow similarité de structures et fonctions
(évolution de gènes/protéines : duplication+modification)

Régions conservées : importance pour la fonction/structure.

EXEMPLE

protéine trypsine : souris (P07146 de SWISS-PROT) et grenouille (P70059 de SWISS-PROT)

```
souris    MSALLILALVGA AVAFPVDDDDKIVGGYTCRESSVPYQVSLNAGYHFCGGSLINDQWVVSAAHC
grenouille MKFLVILVLLGA AVAFEDDD--KIVGGFTCAKNAV PYQVSLNAGYHFCGGSLINSQWVVSAAHC
```

Alignement des deux séquences : représente leur similarité.

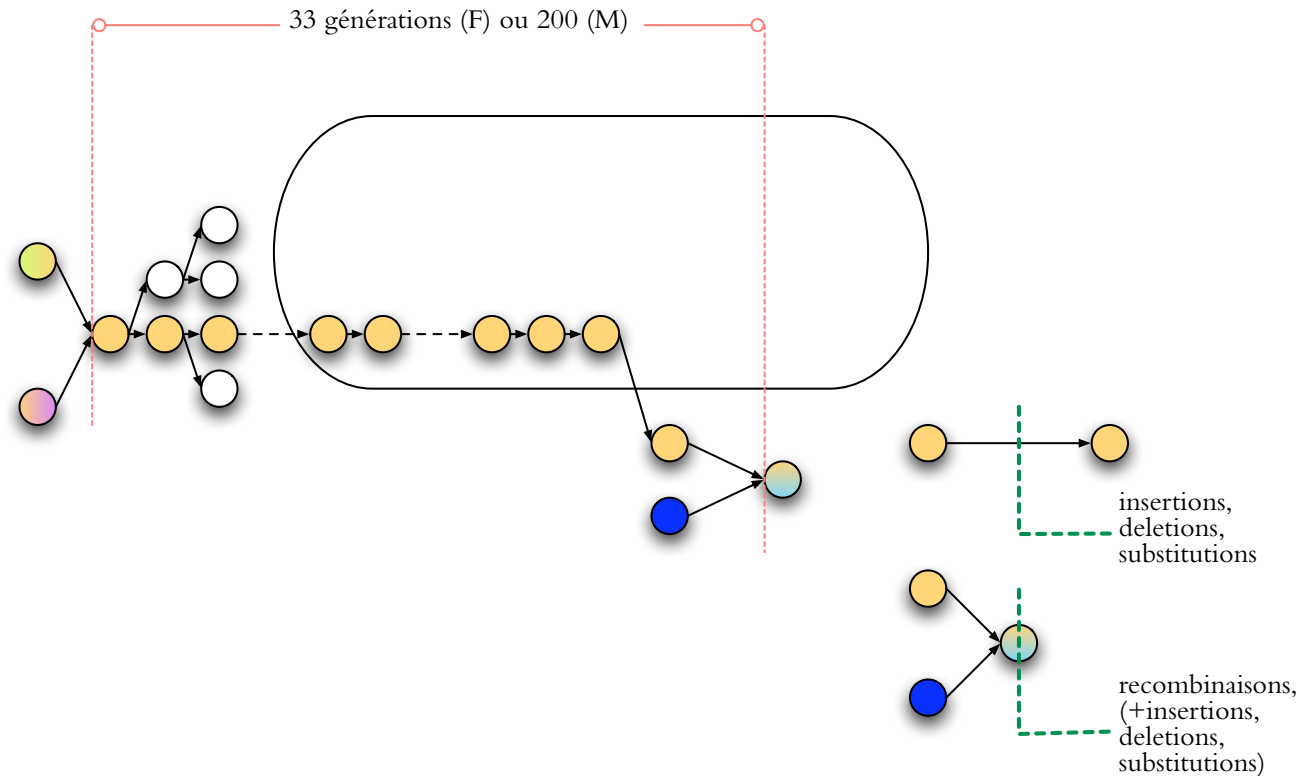
Alphabet d'alignement : $\Sigma \cup \{-\}$.

Appariement : un couple de $\Sigma \cup \{-\}$ (p.e. $(C, -)$, ou (C, G))

Alignement : suite d'appariements.

Relation entre opérations d'éditeurs et alignements (procédure et produit).

MUTATIONS



les mutations s'accumulent pendant les années \rightarrow divergence de séquences d'origin commun

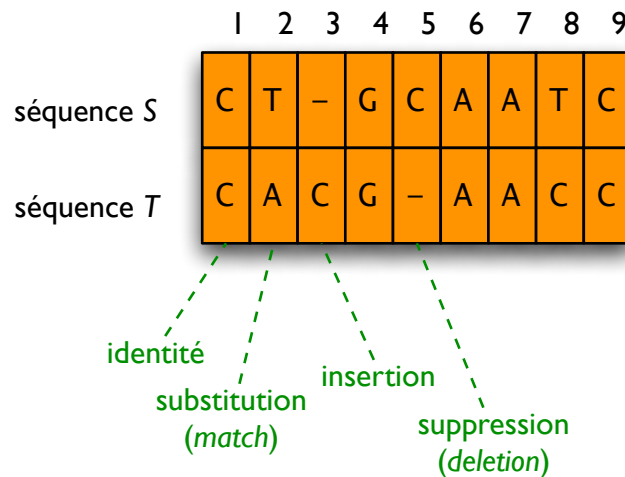
grand divergence \Rightarrow plus de temps et/ou moins de sélection négative

ALIGNEMENT

Déf. *alignement* = vecteur d'appariements

Types d'appariements :

- (a, a) : occurrence (*match*)
- (a, b) : erreur (*mismatch/substitution*)
- $(a, -)$ et $(-, a)$: indel



ALIGNEMENT PONDÉRÉ

Tableau C de pondération d'appariements

	A	C	G	T	—
A	1	-3	-3	-3	-5
C	-3	1	-3	-3	-5
G	-3	-3	1	-3	-5
T	-3	-3	-3	1	-5
—	-5	-5	-5	-5	X

ou :

	A	C	G	T	—
A	91	-114	-31	-123	-300
C	-114	100	-125	-31	-300
G	-31	-125	100	-114	-300
T	-123	-31	-114	91	-300
—	-300	-300	-300	-300	

Score ou valeur de l'alignement : somme des valeurs des appariements

Problème : trouver l'alignement avec le score maximal
 PD pour maximiser la similarité

GÉNÉRATION D'UNE MATRICE DE PONDÉRATION

Modèle probabiliste pour une séquence : caractères aléatoires iid (indépendants et identiquement distribués) : $\mathbb{P}\{S[i] = \sigma\} = \pi_\sigma$

P.e., $\pi_A = \pi_T = 32\%$, $\pi_C = \pi_G = 18\%$: taux de (G + C) à 36%.

Modèle probabiliste pour évolution $S \rightarrow T$: substitutions aléatoires indépendantes : $\mathbb{P}\{T[i] = \sigma' \mid S[i] = \sigma\} = p_{\sigma \rightarrow \sigma'}$ (pas de trous !)

Matrice de substitutions : $\mathbf{M} = \left[p_{\sigma \rightarrow \sigma'} \right]_{\sigma, \sigma' \in \Sigma}$

MATRICE DE PONDÉRATION 2

Probabilité d'un «vrai» alignement : $P_1 = \mathbb{P}\{S, T\}$

Probabilité d'un alignement au hasard : $P_0 = \mathbb{P}\{S\}\mathbb{P}\{T\}$

$$P_1 = \prod_{i=1}^n \mathbb{P}\{S[i], T[i]\} = \prod_{i=1}^n \pi_{S[i]} p_{S[i] \rightarrow T[i]}$$

$$P_0 = \prod_{i=1}^n \mathbb{P}\{S[i]\}\mathbb{P}\{T[i]\} = \prod_{i=1}^n \pi_{S[i]} \pi_{T[i]}$$

MATRICE DE PONDÉRATION 3

Rapport de P_0 et P_1 : **LODS** (logarithmes des chances) $\log \frac{P_1}{P_0}$.

On a

$$\text{LODS} = \sum_{i=1}^n \log \frac{\pi_{S[i]} p_{S[i] \rightarrow T[i]}}{\pi_{S[i]} \pi_{T[i]}} = \sum_{i=1}^n \log \frac{p_{S[i] \rightarrow T[i]}}{\pi_{T[i]}}.$$

Score de substitution $\sigma \rightarrow \sigma'$:

$$\mathbf{C} \begin{bmatrix} \sigma \\ \sigma' \end{bmatrix} = \left\lceil \alpha \log \frac{p_{\sigma \rightarrow \sigma'}}{\pi_{\sigma'}} \right\rceil,$$

où α est un facteur d'échelle (notez qu'on utilise le plafond pour valeurs entières)

MATRICE DE PONDÉRATION 4

Problème : comment estimer π_σ et $p_{\sigma \rightarrow \sigma'}$?

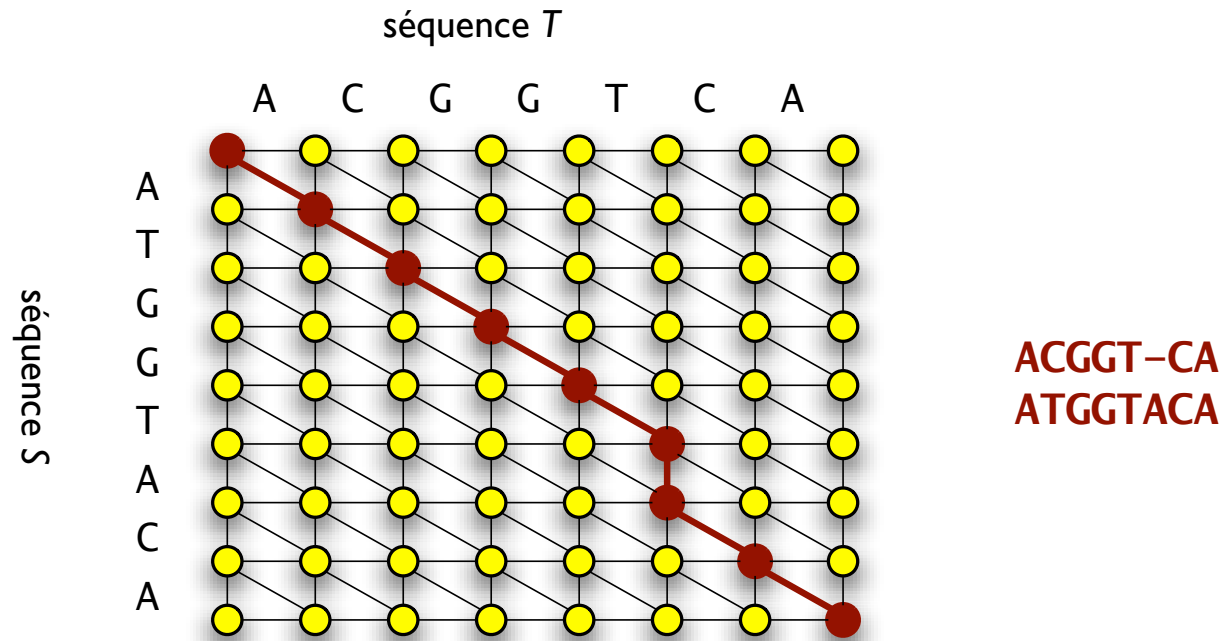
Solution : prendre de «bons alignements»

Exemples : BLOSUM_{nn} (nn=45,62,80 pour nn% d'identité, protéines) ;

PAM_{xxx} (xxx=100,250,..., mesure divergence, protéines) ;

Chiaromonte (pour ADN)

GRAPHE D'ÉDITION



GRAPHE D'ÉDITION 2

pondération des arêtes :

$$\text{poids}(v_{i-1,j-1} \rightarrow v_{i,j}) = \mathbf{C} \begin{bmatrix} S[i] \\ T[j] \end{bmatrix};$$

$$\text{poids}(v_{i-1,j} \rightarrow v_{i,j}) = \mathbf{C} \begin{bmatrix} S[i] \\ - \end{bmatrix};$$

$$\text{poids}(v_{i,j-1} \rightarrow v_{i,j}) = \mathbf{C} \begin{bmatrix} - \\ T[j] \end{bmatrix}.$$

alignement global : trouver le chemin de $v_{0,0}$ à $v_{|S|,|T|}$ avec poids maximal.

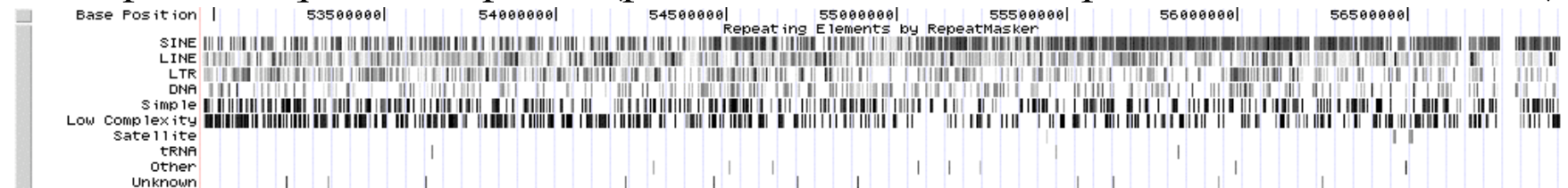
ALIGNEMENT : VARIATION 1

Trouver l'occurrence la plus similaire d'une (courte) séquence dans une autre (longue)

Idée : PD pour les suffixes (ou préfixes) — récurrences similaires

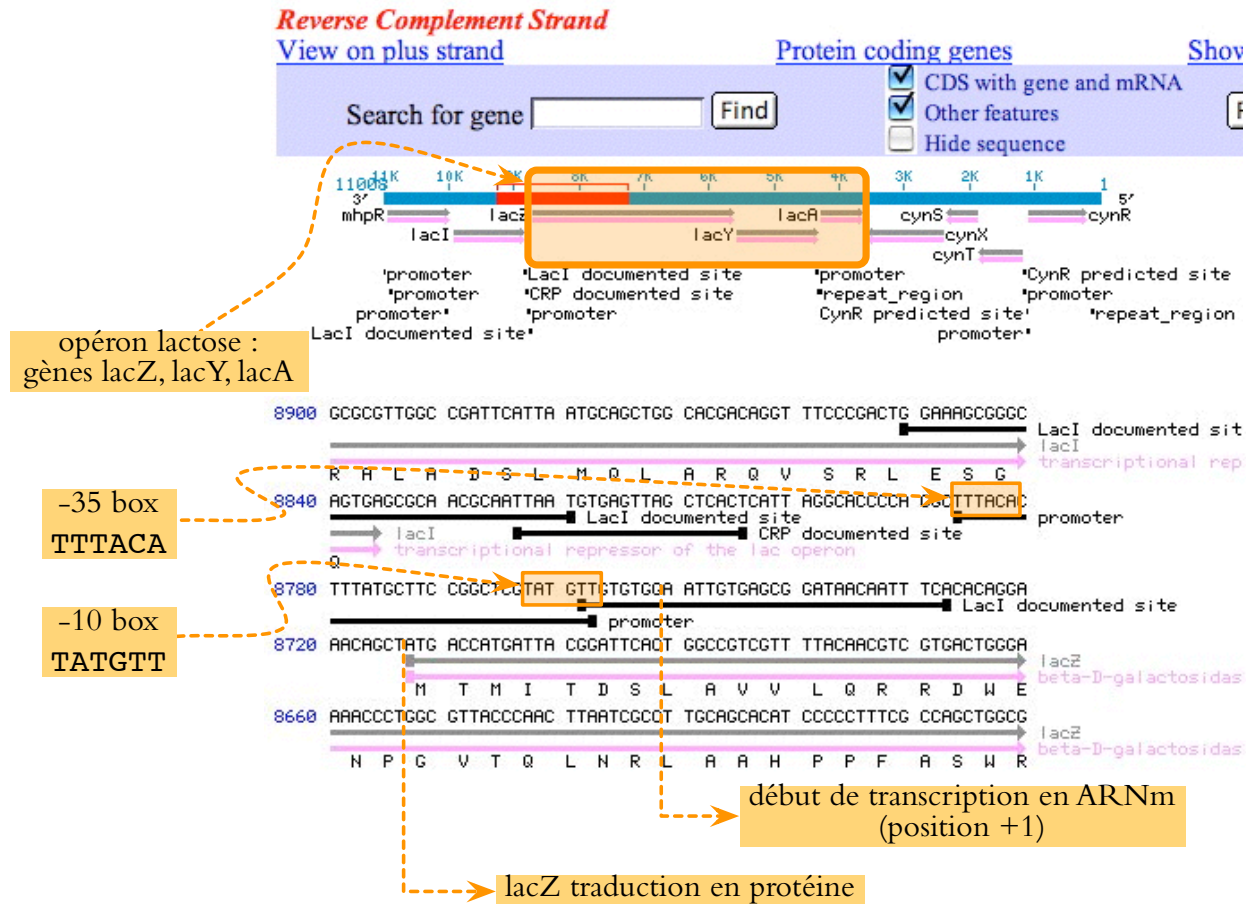
Dans le graphe d'édition : trouver le chemin de $v_{i,0}$ à $v_{i',|T|}$ avec poids maximal pour $i' > i$ quelconques (occurrence de T en S)

Exemple 1 : séquences répétées (p.e., SINE - *Short Interspersed Nuclear Elements*)



Exemple 2 : promoteur sigma A : TTGACA . . . TATAAT [à -35 et -10]

1: AE000141. Escherichia coli ...[gi:1786532]



Brown Genomes, p. 180, Wiley 1999

ALIGNEMENT : VARIATION 2

Trouver la région de similarité maximale entre deux séquences —
alignement local

Dans le graphe d'édition : trouver le chemin de $v_{i,j}$ à $v_{i',j'}$ avec poids maximal pour $i' \geq i, j' \geq j$ quelconques

Exemple : domaines conservés (p.e. homeobox)

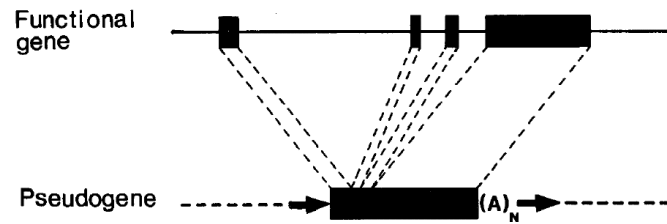
Les noms :

Needleman-Wunsch : problème de l'alignement global

Smith-Waterman : algorithme de l'alignement local

ALIGNEMENT : VARIATION 3

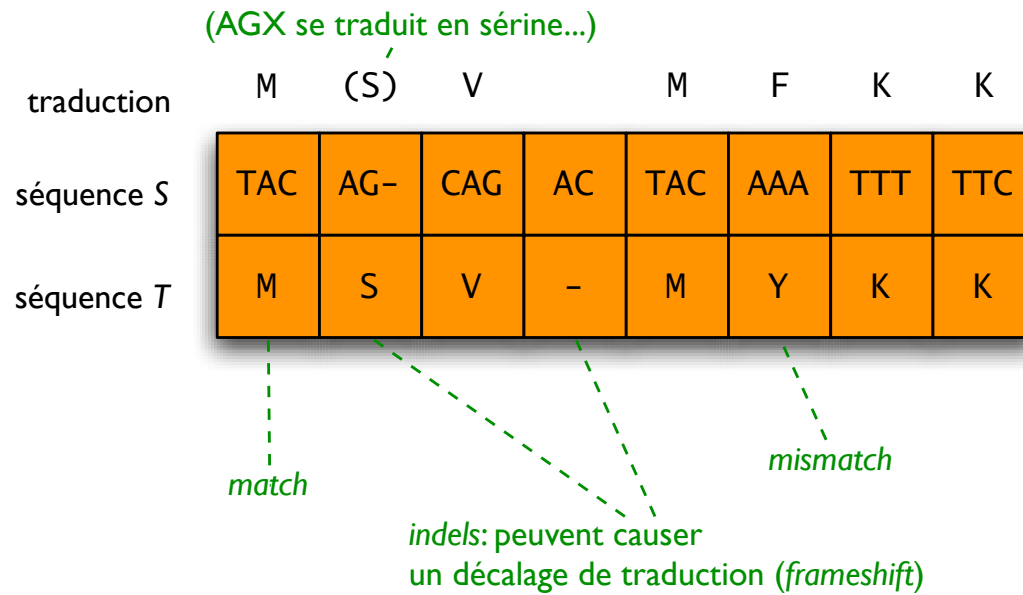
Exemple 1 : Alignement de ADN génomique et ADN codant — pseudogènes



	A	D	Q	T	S	G	D	Q	S	P	L	P		P	C	T	P	T	P	P
SHP2	GCG	GAC	CAG	ACG	AGT	GGA	GAT	CAG	AGC	CCT	CTC	CC--G	CCT	TGT	ACT	CCA	ACG	CCA	CCC	
SHP2-P3	.T.GCT--.A.T.

Vanin *Annu Rev Genet* 19 :253 ; Andersen & al *FASEB J* 18 :8

ALIGNEMENT ADN-PROTÉINE



TROUS

Pénalisation d'un trou par sa taille :

- quelconque : $\delta(\text{longueur})$

- constante

- linéaire : $\delta(\ell) = \delta_{\text{ouvrir}} + \ell\delta_{\text{cont}}$ (avant on avait le cas spécial $\delta_{\text{ouvrir}} = 0$)

Récurrence pour pondération «quelconque» :

$$A(i, j) = \max \left\{ A(i-1, j-1) + \mathbf{C} \begin{bmatrix} S[i] \\ T[j] \end{bmatrix}, \right. \\ \left. \max_{\ell=1, \dots, i} \{A(i-\ell, j) + \delta(\ell)\}, \max_{\ell=1, \dots, j} \{A(i, j-\ell) + \delta(\ell)\} \right\}.$$

Cas de base : $A(i, 0) = \delta(i)$; $A(0, j) = \delta(j)$.

Temps de calcul : $O(nm^2 + mn^2)$ pour $|S| = n$, $|T| = m$.

PONDÉRATION — EXEMPLE

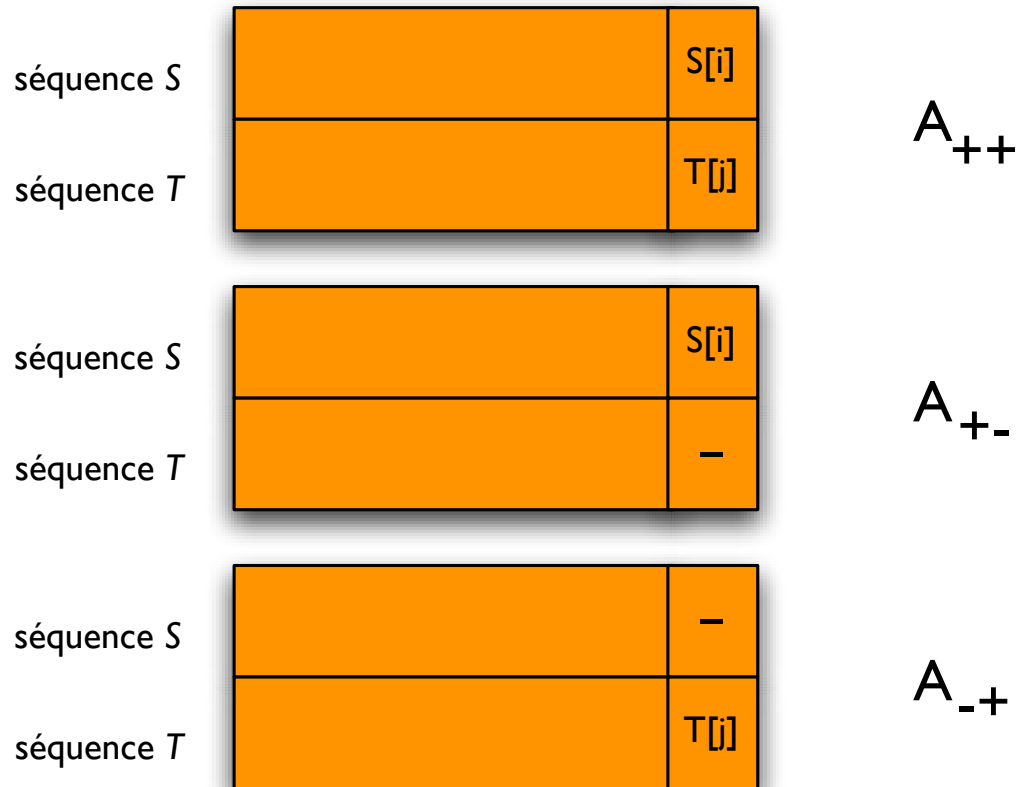
Trous longs (disons $\ell \geq 5$) en S avec pénalisation constante δ_{long} , trous courts avec pénalisation quelconque $\delta(\ell)$; trous en T avec pénalisation simple $\delta'\ell$.

Définir graphe d'alignement avec deux genres de sommets : $v_{i,j}$ pour alignements de préfixes qui finissent pas par trous longs et $t_{i,j}$ pour ceux qui finissent par trous longs.

Arêtes : $v_{i,j} \searrow v_{i+1,j+1}$ ponderée par $\mathbf{C} \begin{bmatrix} S[i] \\ T[j] \end{bmatrix}$, $v_{i,j} \downarrow v_{i+k,j}$ pour $k = 1, \dots, 4$ ponderée par $\delta(k)$, $v_{i,j} \rightarrow v_{i,j+1}$ ponderée par δ' , $u_{i,j} \rightarrow v_{i,j}$ ponderée par 0, $v_{i,j} \downarrow u_{i+5,j}$ ponderée par δ_{long} , $u_{i,j} \downarrow u_{i+1,j}$ ponderée par 0.

PONDÉRATION LINÉAIRE

Réurrences pour $A^{+-}(i, j)$, $A^{-+}(i, j)$ et $A^{++}(i, j)$



PONDÉRATION LINÉAIRE 2

$$A^{++}(i, j) = \mathbf{C} \begin{bmatrix} S[i] \\ T[j] \end{bmatrix} + \max \left\{ A^{++}(i-1, j-1), \right. \\ \left. A^{+-}(i-1, j-1), A^{-+}(i-1, j-1) \right\};$$

$$A^{+-}(i, j) = \max \left\{ A^{+-}(i-1, j) + \delta_1, A^{++}(i-1, j) + \delta_0 \right\};$$

$$A^{-+}(i, j) = \max \left\{ A^{-+}(i, j-1) + \delta_1, A^{++}(i, j-1) + \delta_0 \right\};$$

$$A(i, j) = \max \left\{ A^{++}(i, j), A^{-+}(i, j), A^{+-}(i, j) \right\},$$

où $\delta_1 = \delta_{\text{cont}}$ et $\delta_0 = \delta_{\text{ouvrir}} + \delta_{\text{cont}}$.

[on ignore $A^{+-}(i, j) = A^{-+}(i-1, j) + \delta_1$ ici]

PONDÉRATION LINÉAIRE 3

En fait, on peut éliminer A^{++} :

$$A^{+-}(i, j) = \max \left\{ A^{+-}(i-1, j) + \delta_1, A(i-1, j) + \delta_0 \right\};$$

$$A^{-+}(i, j) = \max \left\{ A^{-+}(i, j-1) + \delta_1, A(i, j-1) + \delta_0 \right\};$$

$$A(i, j) = \max \left\{ \mathbf{C} \begin{bmatrix} S[i] \\ T[j] \end{bmatrix} + A(i-1, j-1), A^{+-}(i, j), A^{-+}(i, j) \right\}.$$

Cas de base : $A(0, 0) = A^{+-}(0, 0) = A^{-+}(0, 0) = 0$; $A(i, 0) = A^{+-}(i, 0) = \delta_0 + (i-1)\delta_1$; $A(0, j) = A^{-+}(0, j) = \delta_0 + (j-1)\delta_1$.

PONDÉRATION LINÉAIRE 4

