

Alignement multiple

Entrée: séquence $S = \{S_1, S_2, \dots, S_m\}$

Sortie: séquences S'_1, S'_2, \dots, S'_m

- toutes de la même longueur
- $\forall i=1 \dots m, S'_i$ est obtenu de S_i par l'insertion de gaps "-"
- critères d'optimisation variable

Critère populaire: "sum-of-pairs" (SP)

somme des paires
On additionne le score entre chaque paire de séquences alignées.

ex: S_1 CAACGAT match +2
 S_2 CAAGGGAT autre -1
 S_3 GAACG gap avec gap 0
 S_4 CACGAT

sol: S'_1 C A A C G - - A T
 S'_2 C A A G G G G A T
 S'_3 G A A C G - - -
 S'_4 C A - C G - - A T

Green annotations: $-1 +2 +2 +2 +2 0 0 -1 -1$ and brackets 9 and 5 .

$$\text{score}(S'_i, S'_j) = \sum_{k=1}^{|S'_i|} M[S'_i[k], S'_j[k]]$$

$$SP(S'_1, S'_2, \dots, S'_m) = \sum_{i=1}^m \sum_{j=i+1}^m \text{score}(S'_i, S'_j)$$

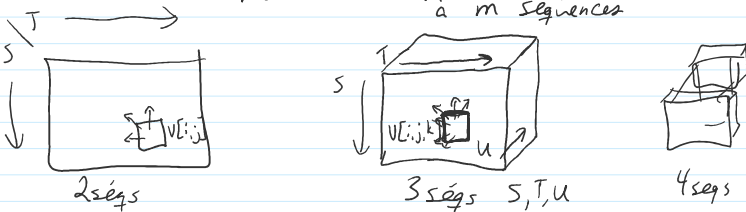
• Trouver les S'_1, S'_2, \dots, S'_m qui minimisent le score SP est NP-complet - sûrement pas d'algo polynomial

• Heuristiques / approximation / algo exact exponentiel

Algo exact exponentiel

Temps $O(n^m \cdot 2^m)$ $n = \text{longueur max des } S_i$
 $m = \text{\# de séquences}$

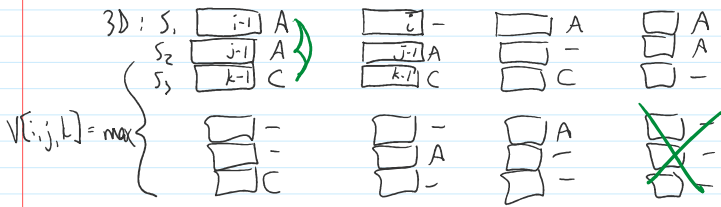
Généralisation de la prog. dynamique pour 2 séquences à m séquences



• $m=3$ (3 séquences S_1, S_2, S_3)

$V[i,j,k] = \text{score optimal d'un alignement multiple entre } S_1[1..i], S_2[1..j], S_3[1..k]$

2D: $\begin{matrix} \square & A \\ \square & A \end{matrix}$ $\begin{matrix} \square & - \\ \square & A \end{matrix}$ $\begin{matrix} \square & A \\ \square & - \end{matrix}$



De façon plus générale

m séquences S_1, S_2, \dots, S_m

$V[i_1, i_2, \dots, i_m] = \text{score opt entre } S_1[1..i_1], \dots, S_m[1..i_m]$

$V[i_1, i_2, \dots, i_m] = \max$

A ou - 2 choix
C ou - 2 choix
...
2^m choix

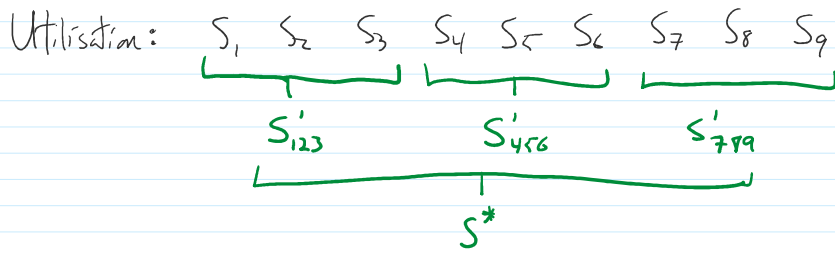
ce max doit énumérer $2^m - 1$ combinaisons de gaps/non-gaps

$$V[i_1, i_2, \dots, i_m] = \max_{(b_1, b_2, \dots, b_m) \in \{0,1\}^m} (V[i_1 - b_1, i_2 - b_2, \dots, i_m - b_m] + \text{score}(b_1, b_2, \dots, b_m))$$

vecteurs de m bits

Complexité: V est une table de dimension $|S_1| \cdot |S_2| \cdot \dots \cdot |S_m|$
 $\in O(n^m)$

Chaque entrée prend un temps $O(2^m)$
 $\Rightarrow O(n^m \cdot 2^m)$

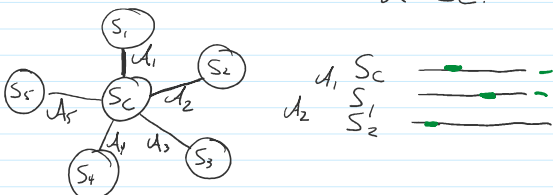


Approximation par centre-étoile

S_1 —
 S_2 —
 \dots —
 S_m —

• Idée: choisir S_c (séquence centrale)
 $c \in \{1, 2, \dots, m\}$

Aligner tous les S_i en fonction de S_c .



Pour choisir S_c , on prend S_i qui maximise

$\sum_{i=1}^m d_i$

😊
 Pour choisir S_c , on prend S_i qui maximise
 $\sum_{j=1}^m \text{score}(S_i, S_j)$ // score = alignement multiple entre 2 séqs

choisir $S_c(S_1, S_2, \dots, S_m)$
 $n = \max |S_i|$
 $m = \# \text{ de séquences}$

m iter
 m iter
 m^2 [pour $i=1..m$
 score = 0
 pour $j=1..m$
 score += alignGlobal(S_i, S_j) $O(n^2)$
 si score > curmax
 curSc = S_i
 curmax = score
 return curSc

Complexité: $O(m^2 \cdot n^2)$

ex: S_1 A C A A T
 S_2 C C C T T
 S_3 A A G A T
 S_4 A G A

On suppose que $S_c = S_1$

A: $S_c = S_1$ A C - A - A T
 S_2 - C - C C T T
 S_3 A - G A - - T
 S_4 A G - A - - -

Ajout de S_2
 Alignement opt $S_c - S_2$
 S_c A C A - A T
 S_2 - C C C T T
 Ajout de S_3
 S_c A C - A - A T
 S_3 A - G A - T
 Ajout de S_4
 S_c A C A A T
 S_4 A G A - -

algo Centre Etoile (S_1, S_2, \dots, S_m)

$O(n^2 m^2)$ $S_c = \text{choisir } S_c(S_1, \dots, S_m)$
 $A = S_c$

// $O(m^2 n^2)$

$O(m)$ iter. pour chaque $S_i \neq S_c$

// Supposition: on stocke chaque align. global calculé dans // choisir S_c

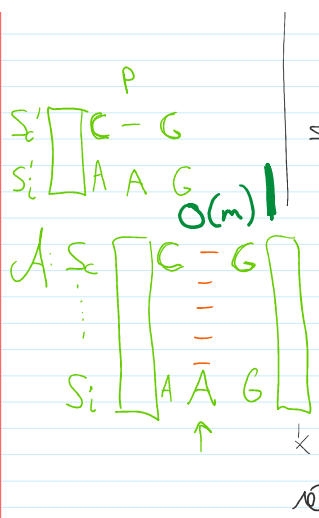
$O(1)$ (S_c', S_i) = alignementGlobal(S_c, S_i)

$O(n)$ iter pour chaque position p de S_c'

si $S_c'[p] \neq "-"$

ajouter le caractère $S_i'[p]$ à la dernière rangée de A , à la position appropriée

P



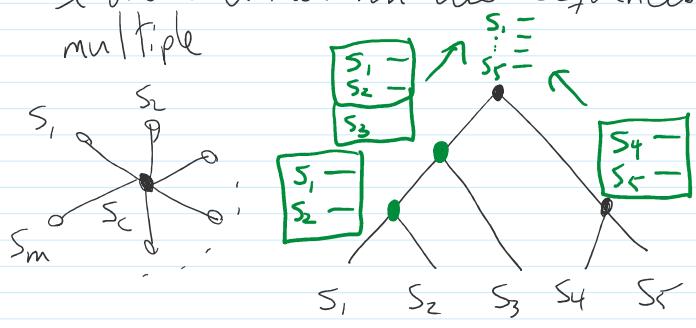
ajouter le caractère $S_i[L_p]$ à la dernière rangée de A , à la position appropriée
 si $S_c[L_p] = "-"$
 ajouter une colonne de gap dans A , entre les positions appropriées, avec $S_i[L_p]$ à la position appropriée

Complexité: $O(m^2 n)$

Si on compte choisir S_c :
 $O(n^2 m^2 + m^2 n) = O(n^2 m^2)$

Clustal

- Utilisation d'un arbre guide pour déterminer l'ordre d'insertion des séquences dans notre align. multiple



// note: noeud \neq feuille
 // noeud = noeud interne + feuilles

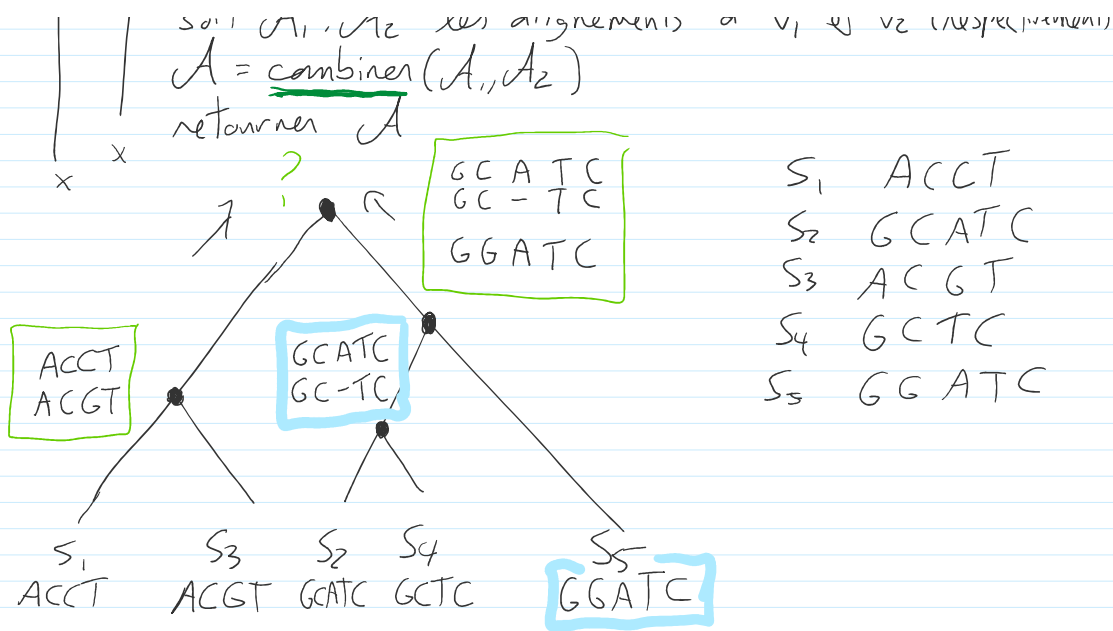
} feuilles = séquences

quand on traite v , tous ses enfants sont traités

clustal (arbre guide T , S_1, S_2, \dots, S_m)

pour chaque noeud v de T dans un parcours post-ordre

si $v =$ feuille,
 | retourner la séquence S_i assignée à v
 sinon
 | soit v_1, v_2 les enfants de v
 | soit A_1, A_2 les alignements à v_1 et v_2 (respectivement)
 | $A =$ combiner (A_1, A_2)



Comment combiner deux alignements multiples?

Entrée: alignements A_1, A_2

Sortie: alignement A avec les rangées de A_1 et A_2

A_1 :

B_1	B_2	B_3	B_4
A	C	C	T
A	C	G	T

On traite B_1, B_2, B_3, B_4
comme des symboles
d'un nouvel alphabet

A_2 :

D_1	D_2	D_3	D_4	D_5
G	C	A	T	C
G	C	-	T	C
G	G	A	T	C

Même chose avec les D_i .

On doit aligner deux séquences

A_1 $B_1 B_2 B_3 B_4$

A_2 $D_1 D_2 D_3 D_4 D_5$

On doit définir une matrice
de scores sur ce nouvel
alphabet

$M[B_i, D_j] = ?$

$M[B_i, "-"] = ?$

$M["-", D_j] = ?$

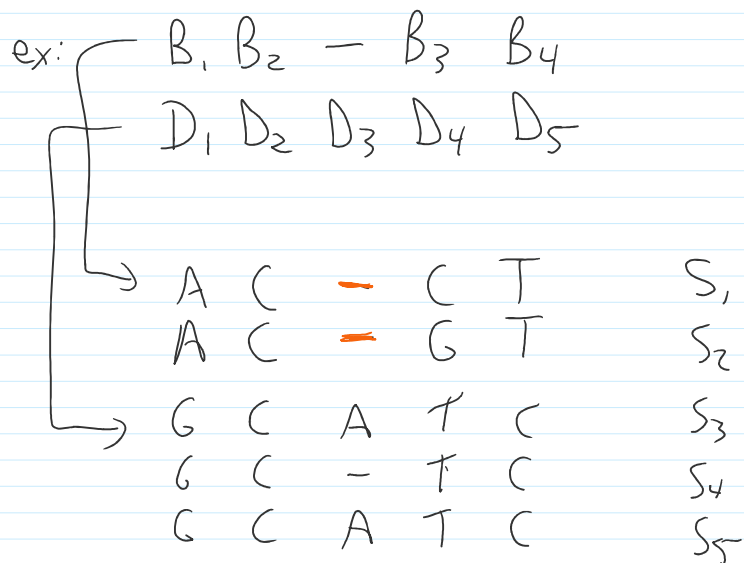
$$M[R, N] = \sum_{i \in B} \sum_{j \in D} M[R, B] \cdot N[D]$$

$$M[B_i, D_j] = \sum_{b=1}^{|B_i|} \sum_{d=1}^{|D_j|} M[B_i[b], D_j[d]]$$

$$\text{ex: } M[B_1, D_1] = M \left[\begin{array}{c} A \\ A \end{array}, \begin{array}{c} G \\ C \\ G \end{array} \right] = \underbrace{M[A, G] + M[A, G] + \dots + M[A, G]}_{2 \times 3 = 6 \text{ fois}}$$

$$M[B_3, D_3] = M \left[\begin{array}{c} C \\ G \end{array}, \begin{array}{c} A \\ - \\ A \end{array} \right] = M[C, A] + M[C, -] + M[C, A] \\ + M[G, A] + M[G, -] + M[G, A]$$

Une fois que $M[B_i, D_j]$ est calculé, $\forall i, j$,
on exécute Needleman-Wunsch avec ce M et
les B_i, D_j



Autres approches: MUSCLE, MAFFT, ...

