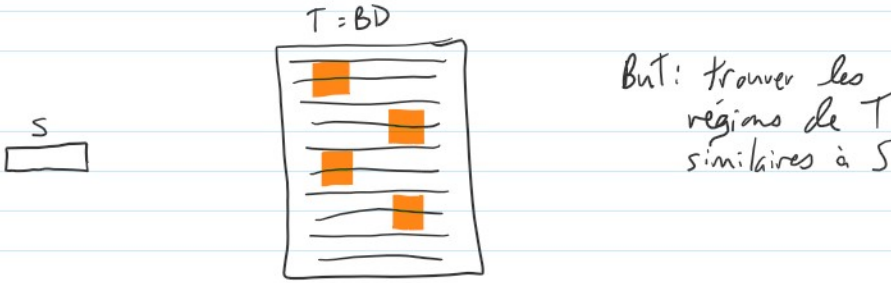


# BLAST

12 septembre 2023 09:13

## Blast Local Alignment Search Tool



But: trouver les régions de  $T$  similaires à  $S$

Alignement local: parfait pour cette tâche mais  $O(|S||T|)$  trop lent  
e.g.  $|T| \sim 100\ 000\ 000\ 000$

Heuristique: algo rapide mais sans garantie de précision

① Prétraitement sur la BD  $T$ : indexation des  $k$ -mers  
 $k$ -mer: séquence de longueur  $k$   $k=3$



ACA: 1  
CAT: 2  
...

$H = \text{dict}(C)$  // clé:  $k$ -mer, val: liste des positions

$O(|T|)$  itér.  $O(k)$  {

pour  $i = 1, \dots, |T| - k + 1$

|  $X = T[i \dots i + k - 1]$

|  $H[X].\text{append}(i)$

x

AAA: 10, 12, 109  
AAC: 1, 4  
:  
i

### ① Recherche de $S$

• Lister tous les  $k$ -mers de  $S$

$k=3$   $S = \text{ACATACCA}$

ACA TAC  
CAT ACC  
ATA CCA

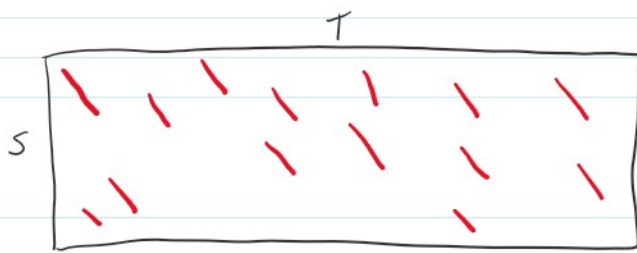
• Lister tous les  $k$ -mers à distance Hamming  $\leq d$  des  $k$ -mers listés (Hamming = #pos différents)

$d=1$  ACA CCA GCA TCA AGA ATA ...  
CAT AAT ...

Pour tous les  $k$ -mers listés, trouver les

Pour tous les k-mers listés, trouver les occurrences dans T. (ut. liser t/)

T = ACA GCA ...

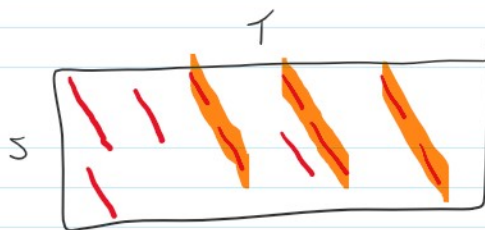


Les occurrences trouvées sont appelées matches

②  $\forall$  match, on étend le match à gauche et à droite maximale, i.e. tant que l'extension améliore le score, jusqu'à un seuil.

S	AAATTC	ACA	TACC	TTTA
T	CGCT	ACAATA	GACC	AAAA
	- - - +	- + + +	- + + +	- - -

③ Joindre les matchs étendus qui sont proches

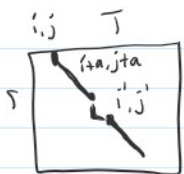


Match  $(S[i, i+a], T[j, j+a])$

$(S[i', i'+b], T[j', j'+b])$

Joindre si  $i' > i+a$   $j' > j+a$   
 et  $i' - (i+a) = j' - (j+a) \leq p$

$p = \text{paramètre}$



④ Retourner les matchs étendus joints

⑤ Exécuter SW sur les matchs trouvés

- ⑤ Exécuter SW sur les matchs trouvés
- ⑥ Trier les matchs par score

E-Value : pertinence d'un score  $s$

$\varepsilon\text{-val}(s)$  = nb espéré d'align. locaux de score  $\geq s$   
 si  $S$  et  $T$  sont des chaînes aléatoire.

Cette espérance tend vers

$$\varepsilon\text{-val}(s) \sim \frac{K \cdot |S| \cdot |T|}{e^{\lambda \cdot s}}$$

$K, \lambda$  dépendent  
de la matrice  
de scores

↓ bas = rare

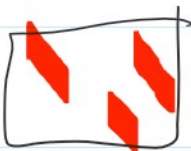
↑ haut = pas rare

$\varepsilon\text{-val}(s) \leq 10^{-5}$  : bon match

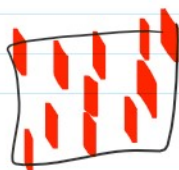
$\leq 10^{-100}$  : très bon match

$\geq 1$  : mauvais

Paramètre  $k$ :  $k$  grand  $\Rightarrow$  - de matches  
 - de précision  
 + rapide



$k$  petit  $\Rightarrow$  + de matches  
 + de précision  
 - rapide



En pratique: nucléotide:  $k=11$

$|S|=4$

En pratique: nucléotide:  $k=11$   $|\Sigma|=4$   
acide aminés:  $k=5$   $|\Sigma|=20$