

Démo 9 – Recherche de structures secondaires de l'ARN

Jean-Eudes Duchesne

1. Mini historique des méthodes de recherche

- i) Algorithmes pour la recherche exacte.
Le but est de trouver un mot P dans une séquence S sans erreurs.
 - Algo naïf, $O(n*m)$
 - Knuth-Morris-Pratt, $O(n)$
 - Boyer-Moore, $O(n)$
- ii) Algorithmes de recherche de motifs simples.
Recherche de motifs tels que des séquences répétées, des séquences miroirs et ainsi de suite. Utilisation d'un arbre de suffixe. Construction de l'arbre se fait en temps $O(n)$. La recherche dans l'arbre peut prendre un temps $O(n)$ ou $O(1)$ tout dépendamment du problème et de la complexité de l'arbre créer.
- iii) Algorithmes de recherche approchée (distance d'édition).
 - Algo naïf, $O(n*m^2)$.
 - Smith-Waterman (alignement semi-local), $O(n*m)$.
 - Algorithmes numériques (Shift-Add, Shift-Or, Shift-And, Myers). Capable d'atteindre $O(n)$ si $p <$ taille d'un mot machine.
- iv) Algorithmes de recherche approchée de structures secondaires.
 - Recherche de molécules biologiques spécifiques, $O(n)$.
 - Recherche générale, $O(sn^2)$ avec grammaires.
 - Modèle d'apprentissage. $O(???)$, mais lent !

2. Classe d'algorithmes de recherche de structures secondaires

Sûr mesure	Générale	Apprentissage
ARNt : tRNAscan, FASrRNA, tRNAscan-se, Bruce, ARAGORN.	RNAMOT, RNAMOTIF, palingol, patscan, Biosmatch.	Covel, ERpin
Intron groupe I : CITRON		
SnoRNA : snoSCAN		

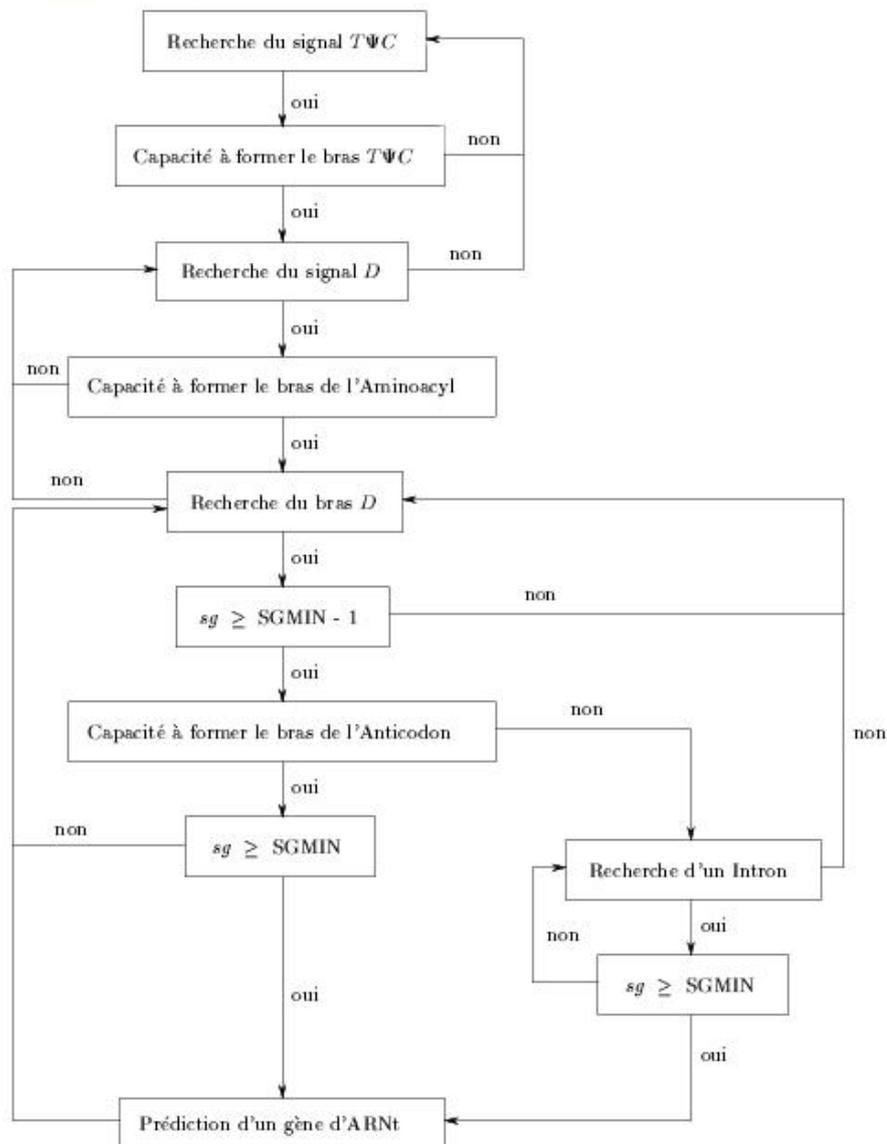
L'approche sûr mesure se veut la plus efficace et la plus rapide des approches. Par contre, celle-ci n'est pas généralisable et très peu flexible. L'idée générale est de rechercher les éléments les plus conservés d'une molécule particulière et d'utiliser ces éléments comme encre pour reconstruire la molécule totale. De façon globale, c'est une approche par satisfaction de contraintes.

L'approche générale utilise plutôt un descripteur pour effectuer sa recherche. Ceci est beaucoup plus flexible puisque c'est l'utilisateur qui entre la séquence à rechercher. Par

contre les résultats dépendent absolument du descripteur. Un bon descripteur donne de bons résultats et inversement pour un mauvais descripteur. Aussi, ces méthodes sont souvent plus lentes puisque des descripteurs avec un grand pouvoir expressif représente un problème de recherche plus complexe que la recherche de motifs conservés. Par contre, ces méthodes permettent à l'utilisateur d'utiliser l'information biologique la plus récente puisque c'est l'utilisateur qui est responsable de mettre à jour le descripteur. Dans le cas de l'approche sûr mesure, l'utilisateur dépend des mises à jours de celui qui maintient le code.

L'approche par apprentissage ressemble à l'approche générale dans le sens où l'algorithme de recherche utilise aussi un descripteur, mais ici l'utilisateur n'est pas responsable de la confection du descripteur. Le programme prends en entrée un ensemble de séquences de la même familles (ex : ARNt de plusieurs espèces différentes) et détermine lui-même les éléments structuraux à rechercher. Dans ce cas, la phase de construction du descripteur est souvent plus longue que la phase de recherche.

3. Exemple sûr mesure : FAStrNA



4. Exemple général : RNAMOT et RNAMOTIF

Descripteur RNAMOT :

H1 s1 H2 s2 H2 s3 H3 s4 H3 s5 H1

H1 3:5 0

H2 4:5 1 AGC:GCU

H3 4:5 1

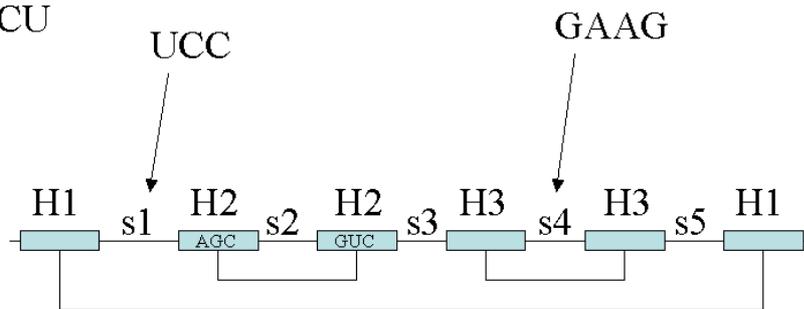
S1 3:6 ucc

S2 5:7

S3 0:3

S4 5:8 gaag

S5 3:5



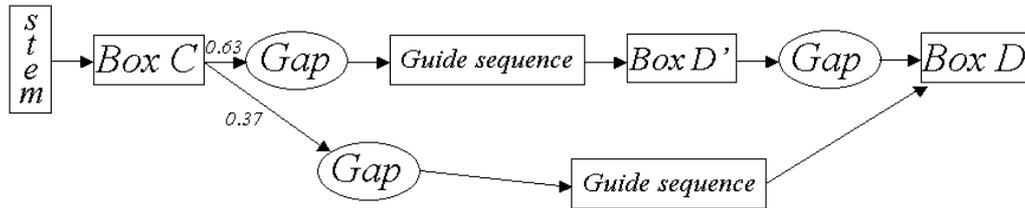
RNATMOTIF est une version améliorée de RNAMOT considéré comme le programme de référence en ce moment dans la recherche de motifs biologiques de façon générale.

RNAMOTIF :

- | | |
|--|---|
| <ol style="list-style-type: none"> 1. Compilation du descripteur (vérification et construction de l'arbre de recherche). 2. Algorithme de recherche (parcours de l'arbre en profondeur d'abord, backtraking). 3. Attribution du score (optionel). | <p>Descripteur <u>rRNA 5s régionIII</u></p> <p><u>parms</u>
<u>wc += gu;</u></p> <p><u>descr</u>
<u>h5(len=2, seq="^AC\$")</u>
<u>h5(len=4, seq="^CYGN\$")</u>
<u>ss(tag='p1', seq="^YCCCATNCCGAAC\$")</u>
<u>h3</u>
<u>ss(len=2)</u>
<u>h3</u></p> |
|--|---|

5. Exemple semi probabiliste : SNOSCAN

SNOSCAN est un peu dans une classe à part dans le sens qu'il utilise les concepts de la recherche par contraintes comme pour les autres programmes de l'approche sûr mesure, mais utilise aussi l'approche probabiliste.



7. Recherche de structures secondaires avec une grammaire

Maintenant, que nous avons un algorithme capable d'aligner une grammaire hors contexte avec une séquence, étendons le problème à des règles de grammaires capables de représenter les hélices biologiques.

Comme nous l'avons précédemment vu, une hélice peut être représentée par les règles suivantes :

$$X \rightarrow n Y n', \quad X \rightarrow n Y, \quad X \rightarrow Y n, \quad X \rightarrow n$$

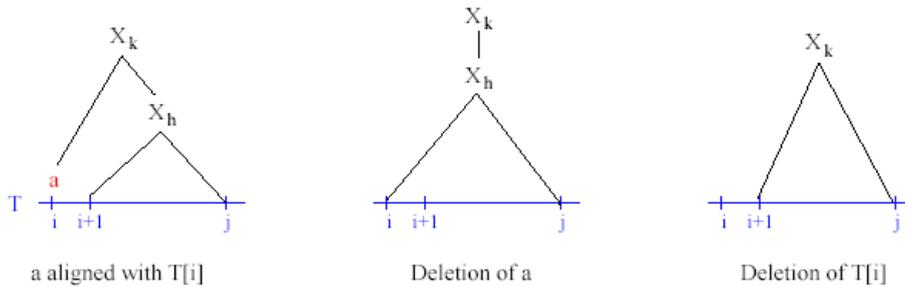
Étant donné une séquence d'ADN, maximisez le score de l'alignement de la grammaire !

- $\delta(a, b)$: Score associé à l'alignement de a et b
- $\delta(a)$: Score associé au indel de a
- $p(a)$: Score d'un appariement correct (a avec a')

Exemple :

$$X_k \rightarrow a X_h$$

$(M_k)_{i,j}$: Max over all X_h of three values



Algorithme :

$$G(X,i,j) = \max \left\{ \begin{array}{l} \max_{X \rightarrow aY} \{ \delta(a, a_i) + G(Y, i+1, j), \\ \delta(a) + G(Y, i, j) \}, \\ \delta(a_i) + G(X, i+1, j), \\ \max_{X \rightarrow Ya} \{ \delta(a, a_i) + G(Y, i, j-1), \\ \delta(a) + G(Y, i, j) \}, \\ \delta(a_i) + G(X, i, j-1), \\ \max_{X \rightarrow aYa'} \{ \delta(a, a_i) + \delta(a', a_j) + p(a) + G(Y, i+1, j-1), \\ \delta(a, a_i) + \delta(a') + G(Y, i+1, j), \\ \delta(a', a_j) + \delta(a) + G(Y, i, j-1), \\ \delta(a) + \delta(a') + G(Y, i, j) \}, \\ \delta(a_i) + G(X, i+1, j), \\ \delta(a_i) + G(X, i, j-1) \end{array} \right. \begin{array}{l} X \rightarrow n Y \\ X \rightarrow Y n \\ X \rightarrow n Y n' \end{array}$$