

Journées PARISTIC
Bordeaux 2005

Bioinformatique / Bioalgorithmique de l'ARN



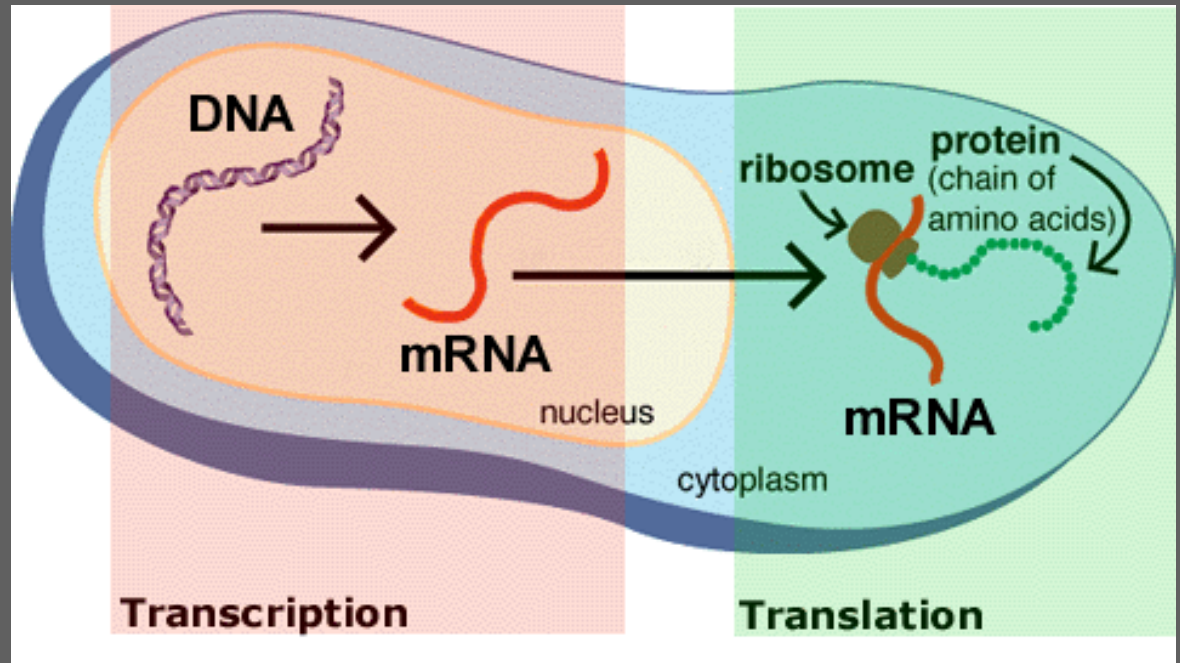
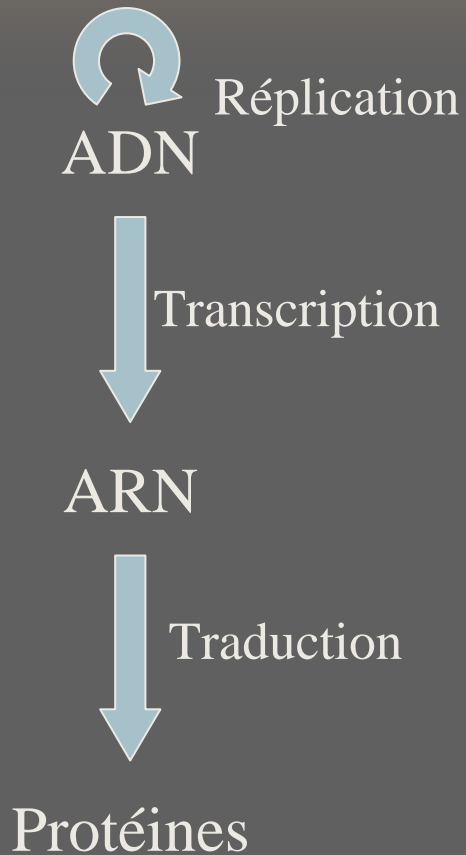
Alain Denise
Bioinformatique
LRI Orsay

UMR CNRS 8623

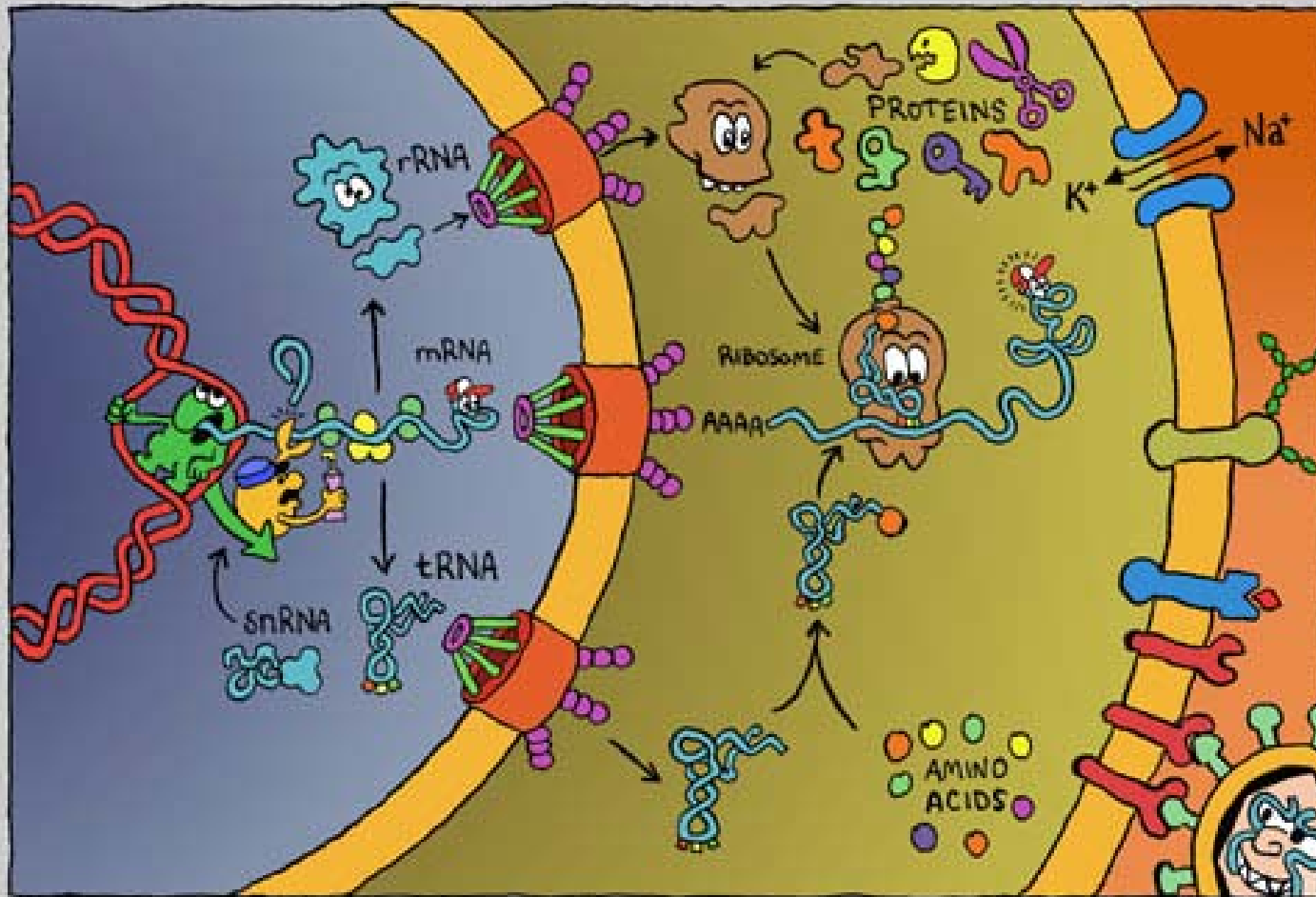
Université Paris-Sud 11



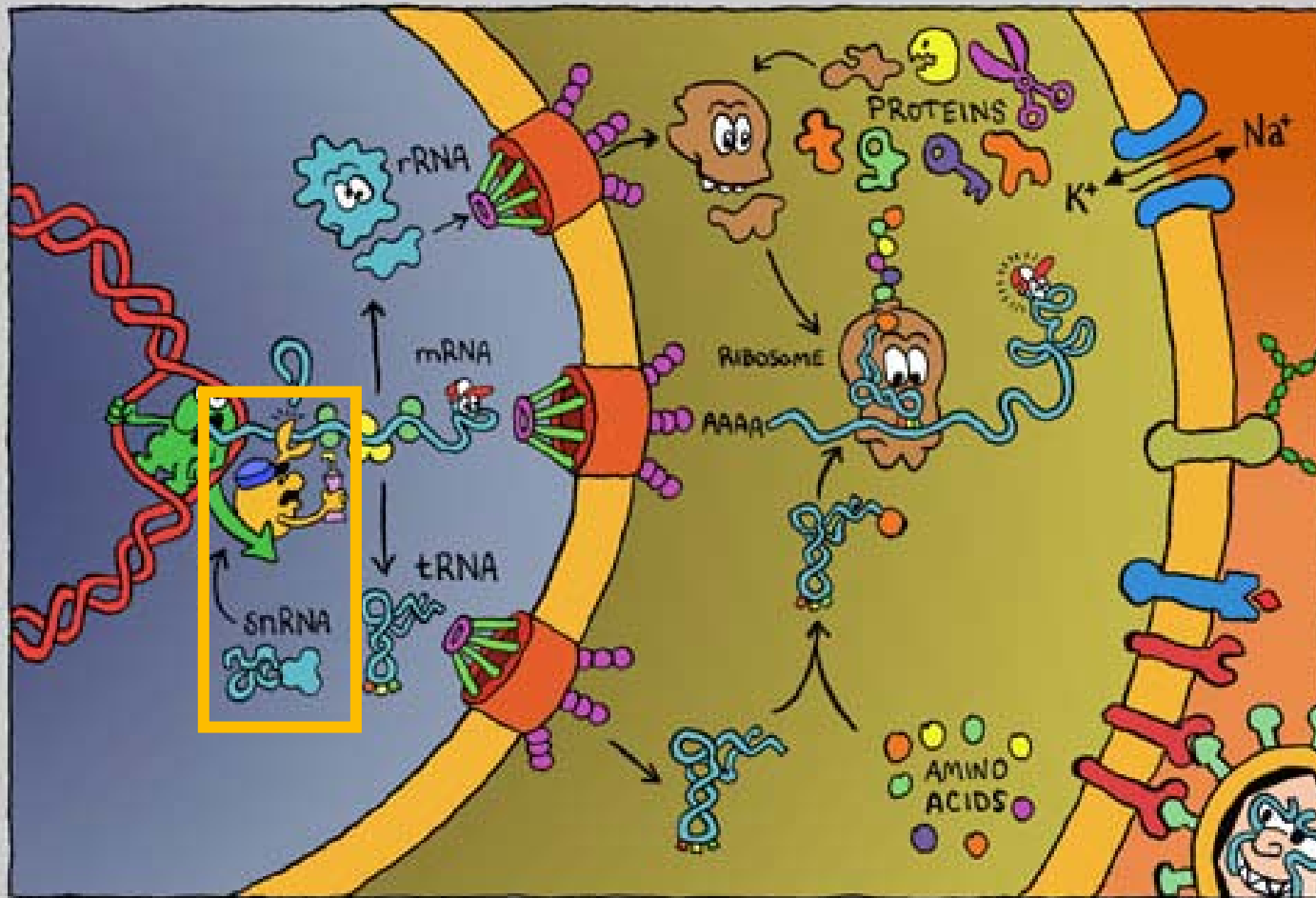
Le dogme central



Les multiples rôles de l'ARN



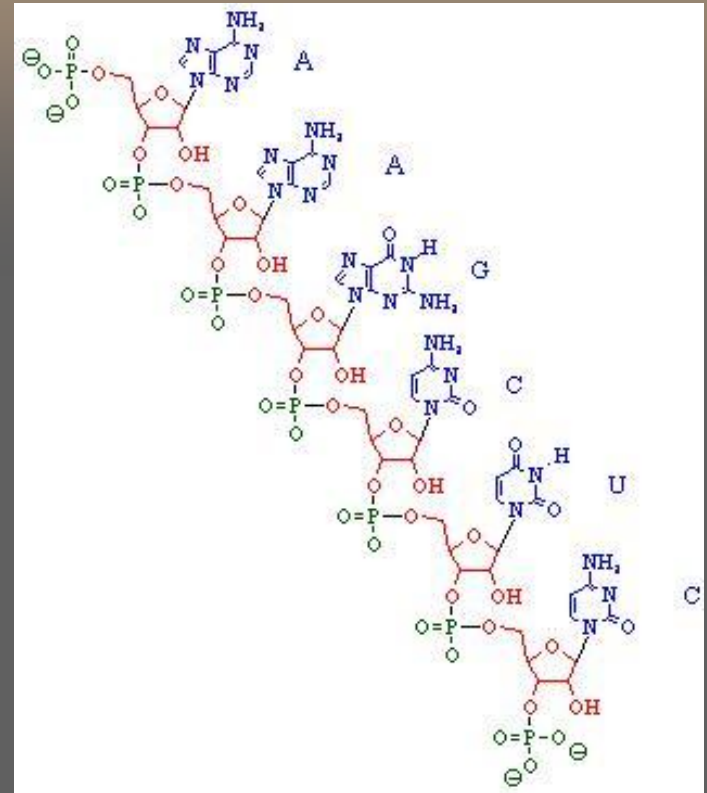
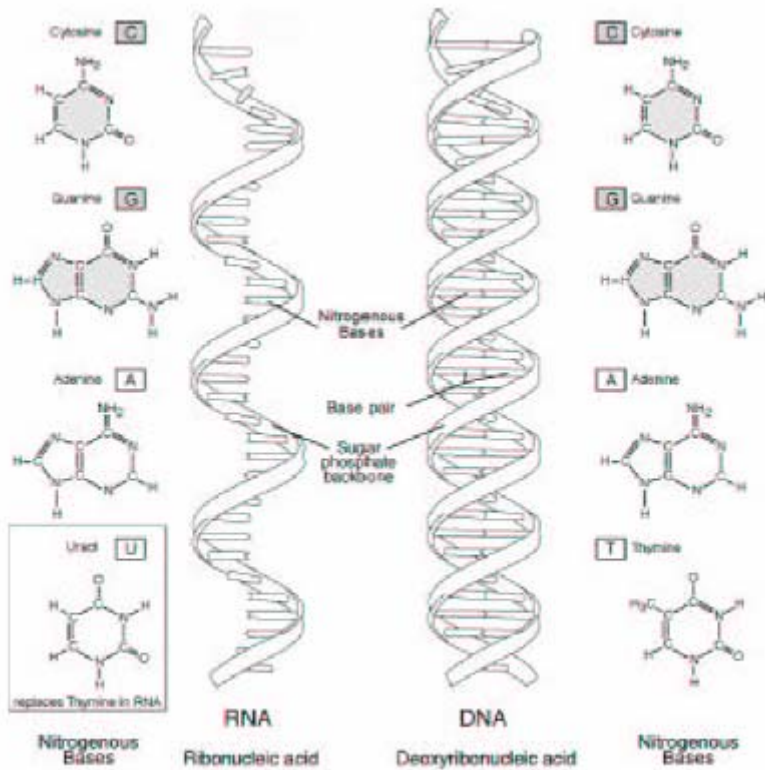
Les multiples rôles de l'ARN



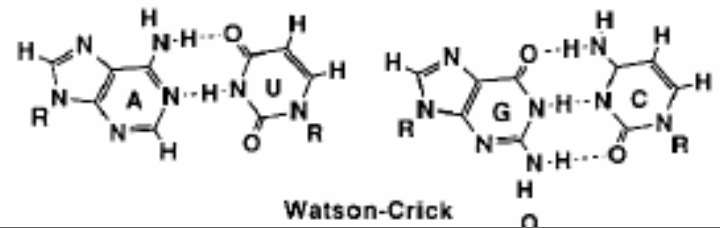
Importance de l'ARN

- ⇒ Présente dans tous les processus cellulaires
- ⇒ La seule molécule qui peut être génome aussi bien que catalyseur
- ⇒ Origine de la vie : le monde à ARN
- ⇒ Cible très fréquente des antibiotiques

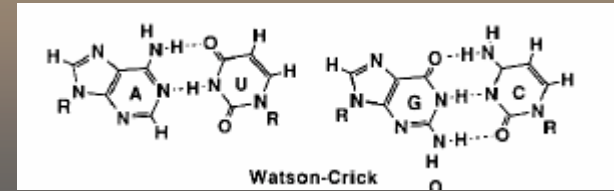
L'ARN



...AAGCUC...



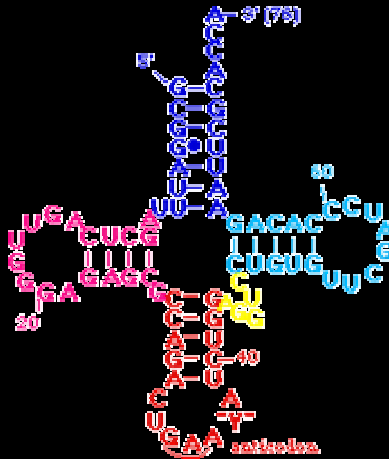
Structure de l'ARN



⇒ Structure primaire

CGCGAUUUAGCUCAGUUGGGAGAGCGCCAGACUGAAUAUCUGGAGGUC CUGUGUUCGAUCCACAGAAUUCGCACCA

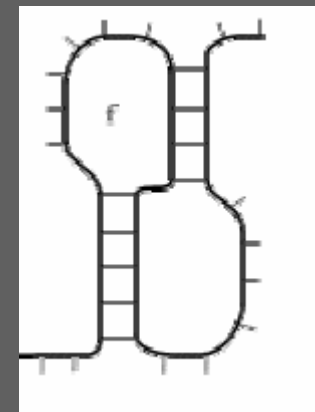
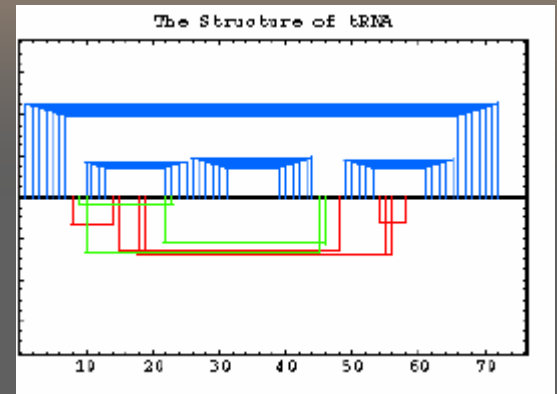
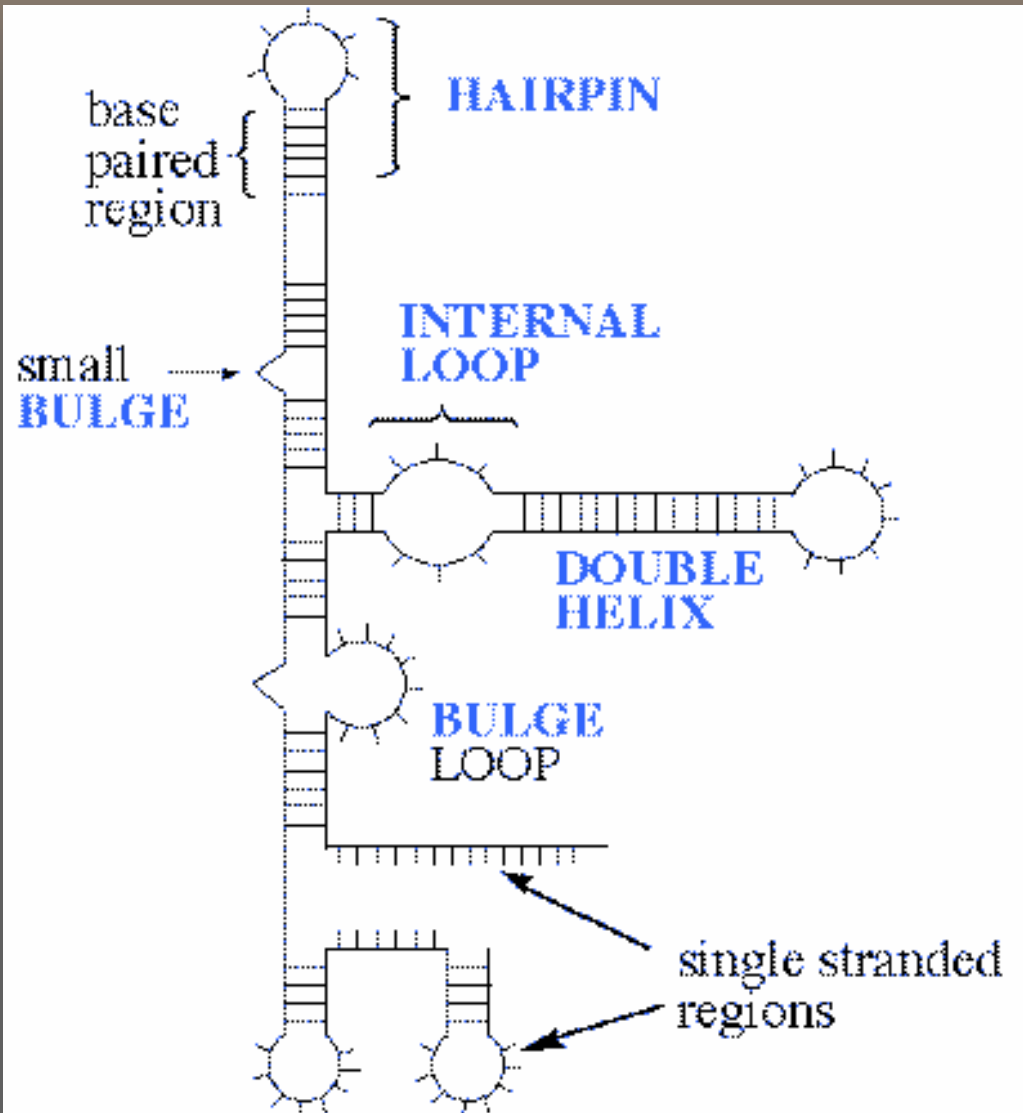
⇒ Structure secondaire



⇒ Structure tertiaire



Structures secondaires

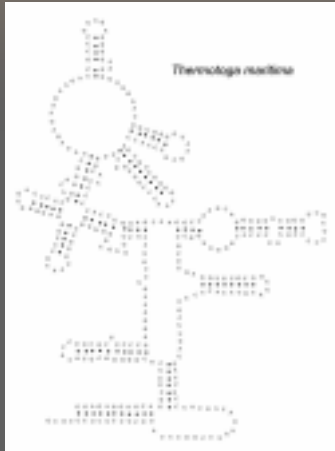


Pseudo-nœud

« Bio-Algorithmique » de l'ARN

- ⇒ Prédiction de structure en fonction de la séquence
- ⇒ Détection de motifs structurels dans une séquence
- ⇒ Comparaison de deux ou plusieurs structures
- ⇒ Détermination d'une séquence en fonction de la structure
- ⇒ Recherche de sous-structures communes à deux ou plusieurs structures

Pourquoi comparer ?



- Phylogénie
- Prédiction de structures
- Recherche de motifs communs



Comparer = Calculer une **distance** (ou un score)

Edition et alignement deux à deux

On se donne un ensemble « d'opérations atomiques », chacune ayant un score (ou un coût).

Données : deux structures.

- **Edition** : trouver la suite d'opérations de score maximal (ou de coût minimal) permettant de transformer une structure en l'autre.
- **Alignement** : trouver une « sur-structure » commune aux deux structures telle que la somme des scores d'édition de chacune des structures à la sur-structure soit maximale (ou que la somme des coûts soit minimale).

Comparaison de 2 séquences

Deux séquences $v = v_1v_2\dots v_n$ et $w = w_1w_2\dots w_m$

Opérations d'édition :

- $\text{ins}(x,i)$
- $\text{suppr}(x,i)$
- $\text{subs}(x,y,i)$

CHAT - $\text{suppr}(C,1) \rightarrow$ HAT - $\text{subs}(H,R,1) \rightarrow$ RAT

(Pour les séquences : édition \sim alignement)

Comparaison de 2 séquences

$$V = v_1 v_2 \dots v_n$$

$$W = w_1 w_2 \dots w_m$$

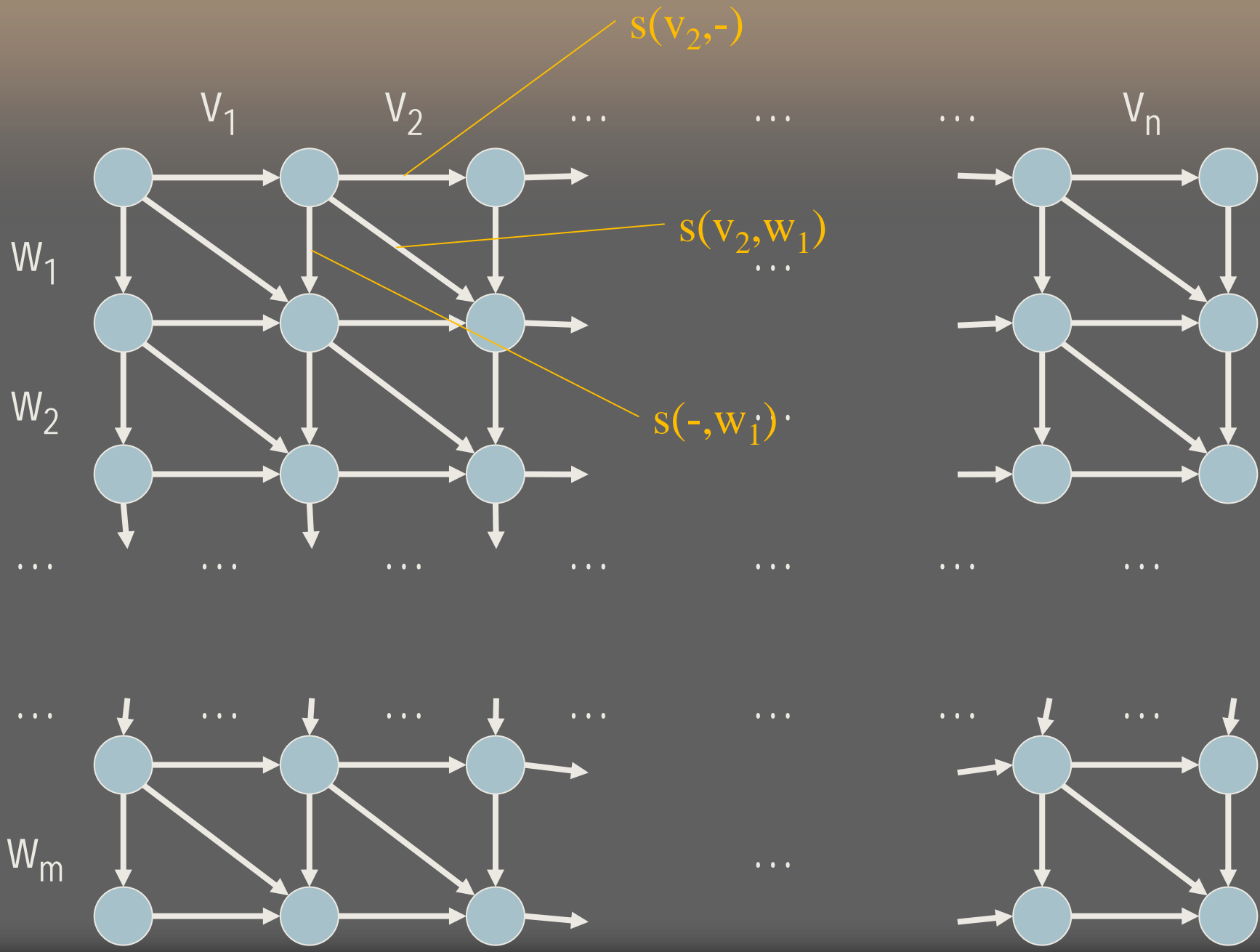
$s(x,y)$: coût de substitution de x en y

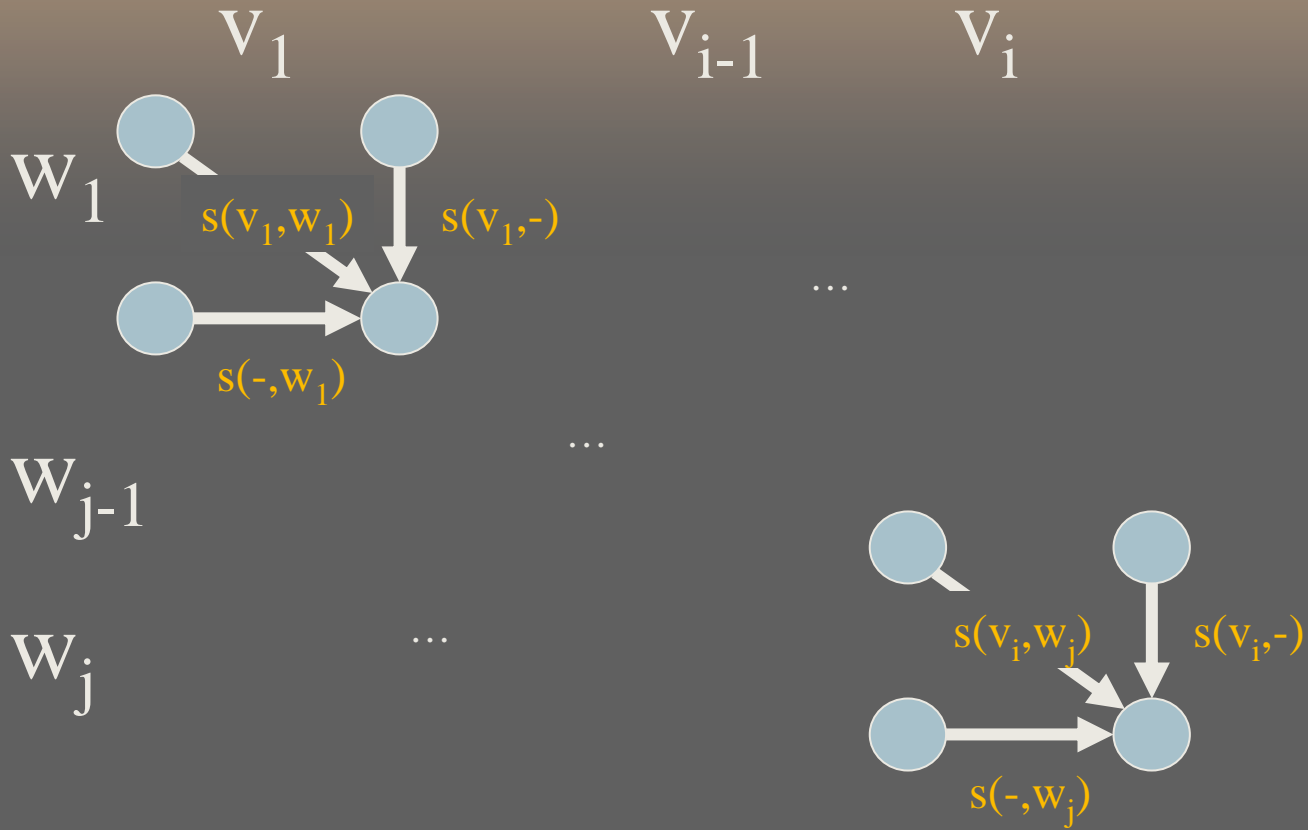
$s(x,-)$: coût de suppression de x

$s(-,y)$: coût d'insertion de y

$D(v,w)$: distance d'édition de v et w

$$D(v_1 \dots v_i, w_1 \dots w_j) = \text{Min} \left\{ \begin{array}{l} D(v_1 \dots v_{i-1}, w_1 \dots w_{j-1}) + s(v_i, w_j) \\ D(v_1 \dots v_{i-1}, w_1 \dots w_j) + s(v_i, -) \\ D(v_1 \dots v_i, w_1 \dots w_{j-1}) + s(-, w_j) \end{array} \right\}$$





$$D(v_1 \dots v_i, w_1 \dots w_j) = \text{Min} \{$$

$$D(v_1 \dots v_{i-1}, w_1 \dots w_{j-1}) + s(v_i, w_j)$$

$$D(v_1 \dots v_{i-1}, w_1 \dots w_j) + s(v_i, -)$$

$$D(v_1 \dots v_i, w_1 \dots w_{j-1}) + s(-, w_j)$$

}

Opérations de séquences arc-annotées

⇒ Opérations sur les bases :

- ▣ Suppression / Insertion
- ▣ Substitution (ou conservation)

⇒ Opérations sur les arcs :

▣ Suppression / Insertion :

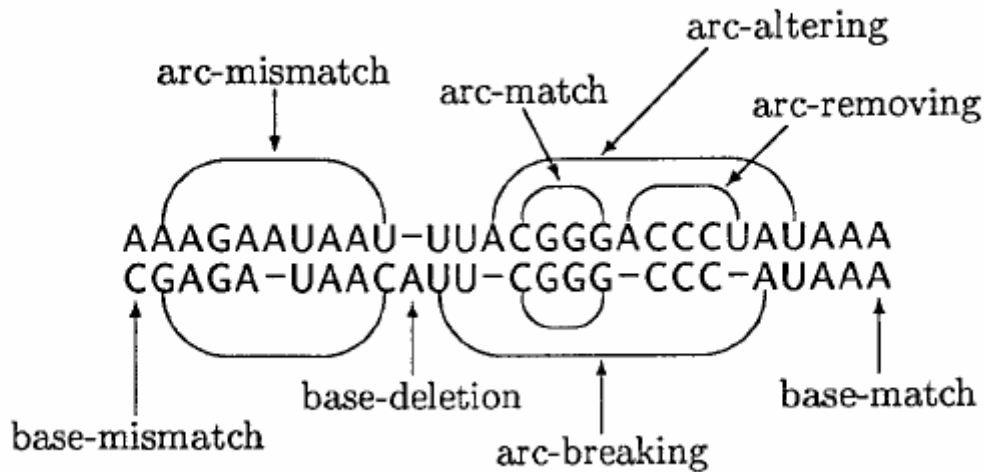
▣ Cassure / :

▣ Altération / :

▣ Substitution :



Edition de séquences arc-annotées

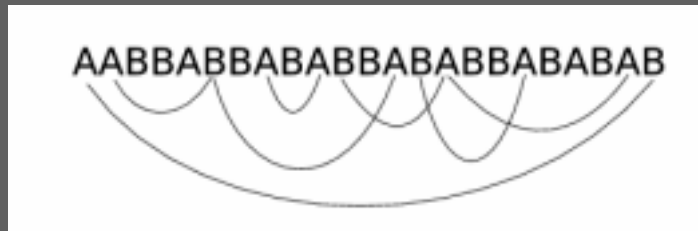


(Jiang, Lin, Ma, Zhang 2002)

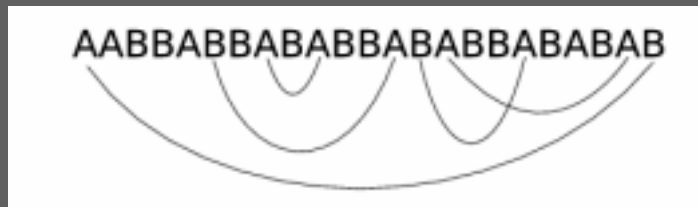
Complexité de l'édition

Types de séquences arc-annotées

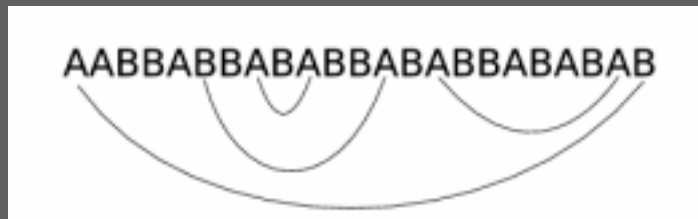
⇒ Générale



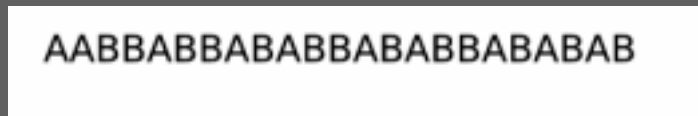
⇒ Croisée









⇒ Imbriquée



⇒ Sans arcs



Complexité de l'édition

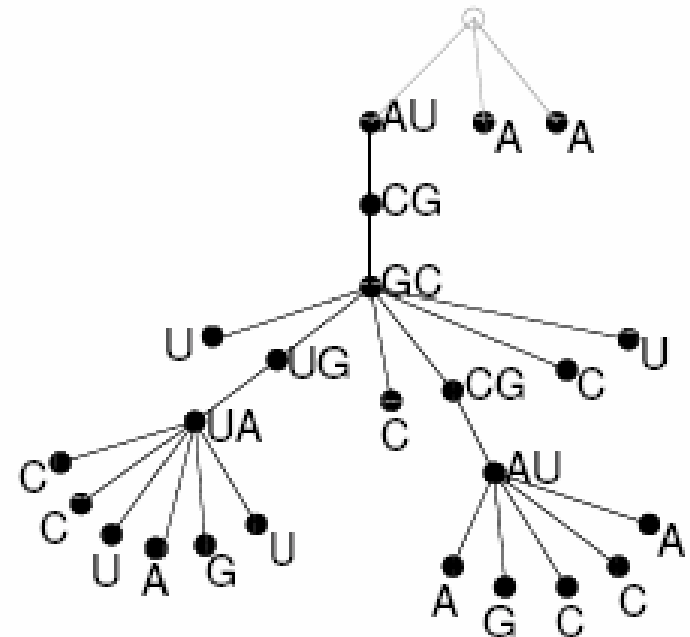
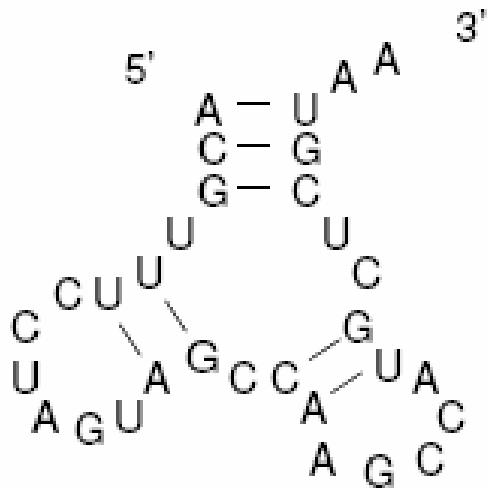
	Générale	Croisée	Imbriquée	Sans arcs
Générale	NP-complet			
Croisée		NP-complet		
Imbriquée			NP-complet	$O(nm^3)$
Sans arcs				$O(nm / \log n)$

Si $2 \times \text{Score}(\text{Altération d'arc}) = \text{Score}(\text{Cassure}) + \text{Score}(\text{Suppression})$, alors
algorithme en $O(n^3m)$ pour $\text{Edit}(\text{croisée}, \text{imbriquée})$ et $\text{Edit}(\text{imbriquée}, \text{imbriquée})$

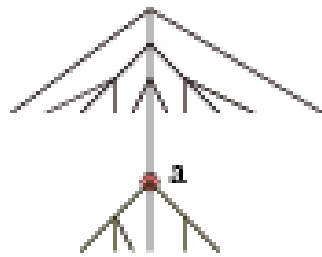
- Jiang, Lin, Ma, Zhang 2002
- **Blin, Fertin, Rusu, Sinoquet 2003**
- Crochemore, Landau, Ziv-Ukelson 2000

Le cas « imbriqué-imbriqué »

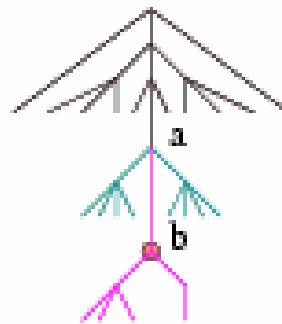
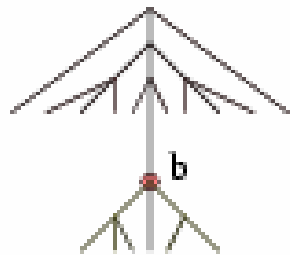
- Structures secondaires
- Comparaison d'arbres



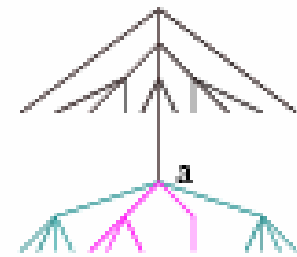
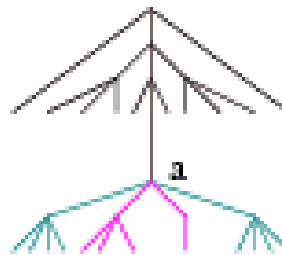
Opérations d'édition



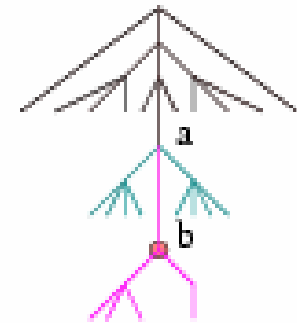
change(a → b)



delete(b)



insert(b)



Algorithme d'édition

Zhang, Shasha 1989

$$\begin{aligned}
 & Tscore(\begin{array}{c} \bullet \\ \swarrow \quad \searrow \\ \blacktriangle \quad \blacktriangle \end{array} , \begin{array}{c} \bullet \\ \swarrow \quad \searrow \\ \blacktriangle \quad \blacktriangle \quad \blacktriangle \end{array}) \\
 = \text{Max} & \left\{ \begin{array}{l} Fscore(\blacktriangle \quad \blacktriangle , \blacktriangle \quad \blacktriangle \quad \blacktriangle) + \text{Change}(\bullet , \bullet), \\ Fscore(\blacktriangle \quad \blacktriangle , \begin{array}{c} \bullet \\ \swarrow \quad \searrow \\ \blacktriangle \quad \blacktriangle \quad \blacktriangle \end{array}) + \text{Delete}(\bullet), \\ Fscore(\begin{array}{c} \bullet \\ \swarrow \quad \searrow \\ \blacktriangle \quad \blacktriangle \end{array} , \blacktriangle \quad \blacktriangle \quad \blacktriangle) + \text{Insert}(\bullet) \end{array} \right\}
 \end{aligned}$$

$$\begin{aligned}
 & Fscore(\begin{array}{c} \blacktriangle \\ \swarrow \quad \searrow \\ \blacktriangle \quad \blacktriangle \end{array} , \begin{array}{c} \bullet \\ \swarrow \quad \searrow \\ \blacktriangle \quad \blacktriangle \quad \blacktriangle \end{array}) \\
 = \text{Max} & \left\{ \begin{array}{l} Fscore(\blacktriangle , \blacktriangle \quad \blacktriangle) + Tscore(\begin{array}{c} \bullet \\ \swarrow \quad \searrow \\ \blacktriangle \quad \blacktriangle \end{array} , \begin{array}{c} \bullet \\ \swarrow \quad \searrow \\ \blacktriangle \quad \blacktriangle \quad \blacktriangle \end{array}), \\ Fscore(\blacktriangle \quad \blacktriangle \quad \blacktriangle , \begin{array}{c} \bullet \\ \swarrow \quad \searrow \\ \blacktriangle \quad \blacktriangle \quad \blacktriangle \end{array}) + \text{Delete}(\bullet), \\ Fscore(\blacktriangle \quad \begin{array}{c} \bullet \\ \swarrow \quad \searrow \\ \blacktriangle \quad \blacktriangle \end{array} , \blacktriangle \quad \blacktriangle \quad \blacktriangle) + \text{Insert}(\bullet) \end{array} \right\}
 \end{aligned}$$

Opérations d'édition

⇒ Opérations sur les bases :

- Suppression / Insertion

- Substitution

⇒ Opérations sur les arcs :

- Suppression / Insertion :

- Cassure / :

- Altération / :

- Substitution :



Opérations d'édition : manques

⇒ Opérations sur les bases :

▣ Suppression / Insertion

▣ Substitution

⇒ Opérations sur les arcs :

▣ Suppression / Insertion :

~~▣ Cassure / :~~

~~▣ Altération / :~~

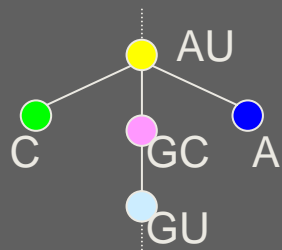
▣ Substitution :



Opérations d'édition : problème

A-U
U-A
G-C
C-U

A-U
C A
G-C
C-U



Delete(●)

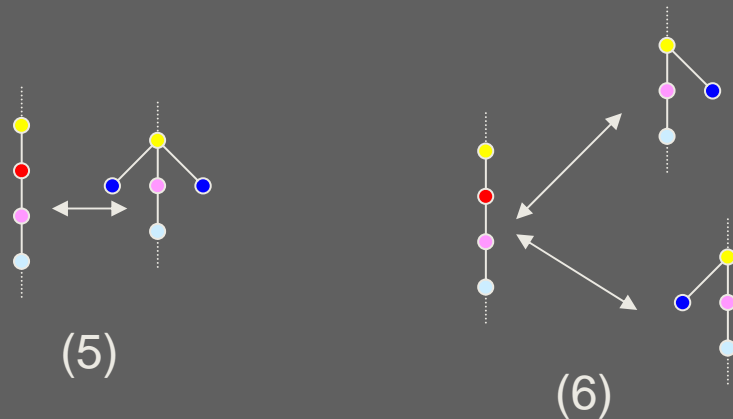
Insert(●)

Insert(●)

} 3 opérations au lieu d'une !

Opérations d'édition : ajouts

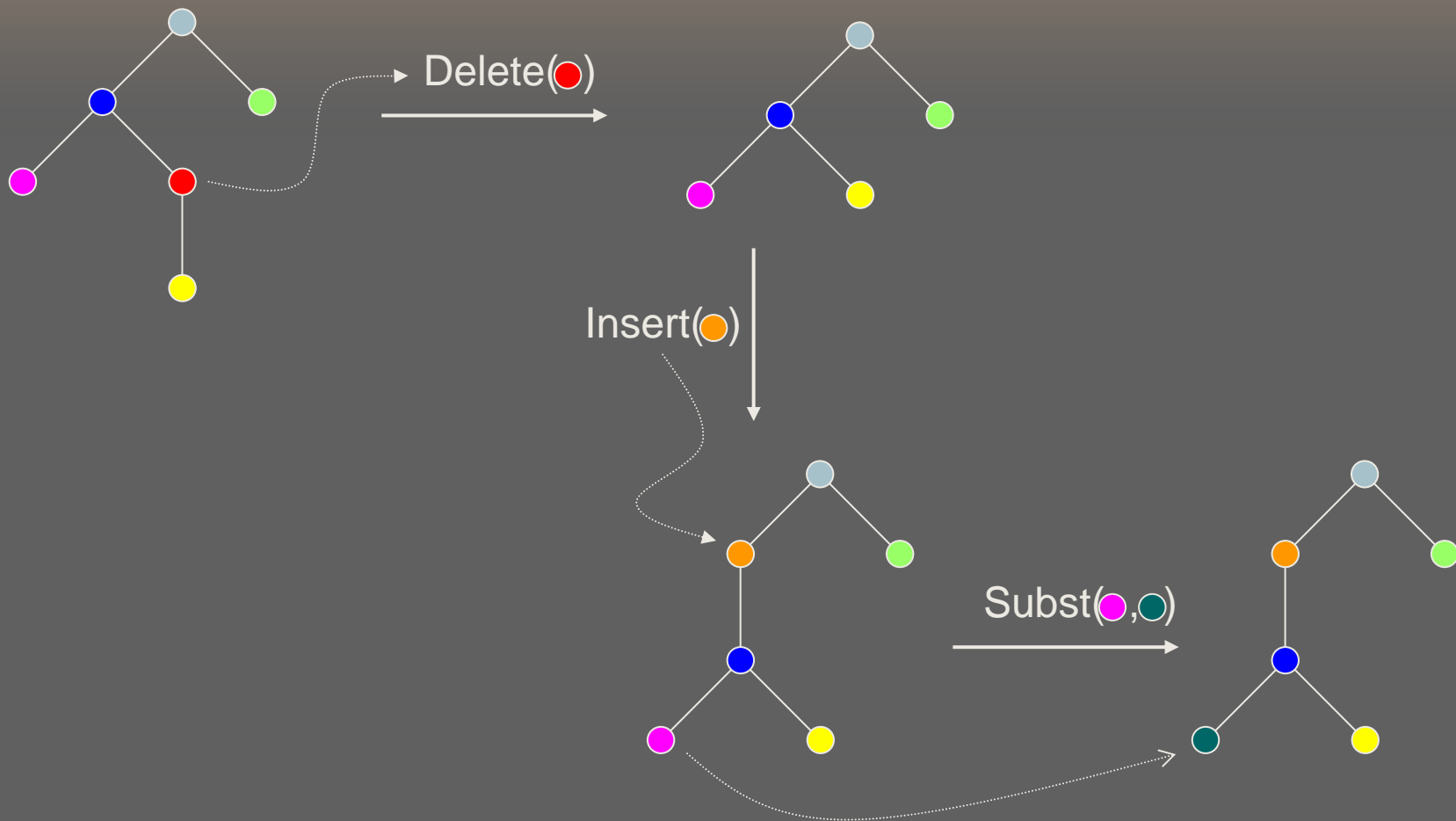
- ⇒ Suppression et insertion d'une base
- ⇒ Substitution de bases
- ⇒ Suppression et insertion d'une paire de bases
- ⇒ Substitution de paires de bases
- ⇒ Appariement et désappariement (5)
- ⇒ Suppression et insertion d'une base dans une paire de bases (6)



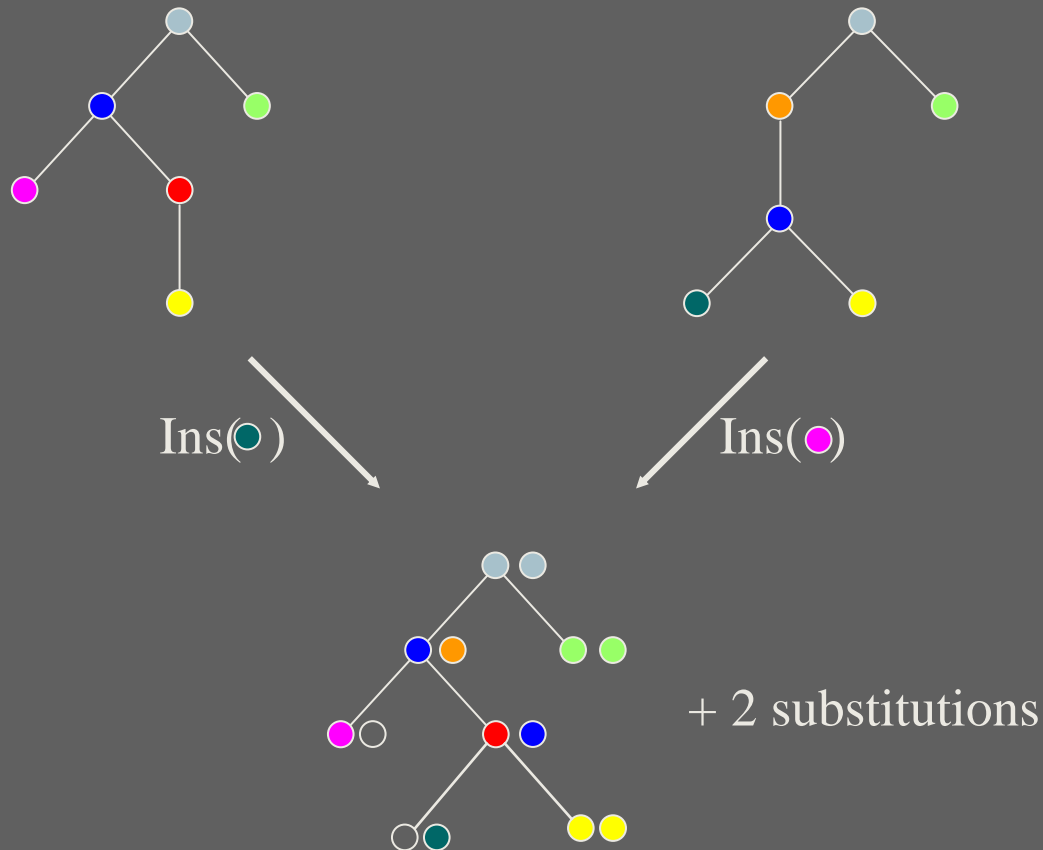
Edition et alignement d'arbres

	Opérations « arbres »	Opérations « ARN »
Edition	$O(n^3 \log n)$ [Zhang-Shasha 1989, Klein 1998]	NP-complet [Blin, Fertin, Sinoquet, Rusu 2003]
Alignement	$O(n^4)$ [Jiang, Wang, Zhang 1995]	?

Edition d'arbres \neq Alignement



Edition d'arbres \neq Alignement



Qu'est-ce que l'alignement dans notre cas?



Edition



Alignement



Un algorithme d'alignement (1/4)

Herrbach, AD, Dulucq, Touzet 200

$$\text{Score}(\text{▲▲▲}, \text{▲▲▲}) = \text{Max}$$

$$\text{Score}(\emptyset, \text{▲▲▲}) + \text{DelB}(\bullet)$$

si \bullet est une base

$$\text{Score}(\text{▲▲▲}, \emptyset) + \text{InsB}(\circ)$$

si \circ est une base

$$\text{Score}(\emptyset, \emptyset) + \text{SubB}(\bullet, \circ)$$

si \bullet et \circ sont des bases

$$\text{Score}(\text{▲▲▲}, \text{▲▲▲}) + \text{DelP}(\bullet)$$

si \bullet est une paire

$$\text{Score}(\text{▲▲▲}, \text{▲▲▲}) + \text{InsP}(\circ)$$

si \circ est une paire

$$\text{Score}(\text{▲▲▲}, \text{▲▲▲}) + \text{SubP}(\bullet, \circ)$$

si \bullet et \circ sont des paires

$$\text{Score}(\text{▲▲▲}, \emptyset) + \text{Del5}(\bullet, \circ)$$

si \bullet une paire et \circ une base

$$\text{Score}(\text{▲▲▲}, \emptyset) + \text{Del3}(\bullet, \circ)$$

si \bullet une paire et \circ une base

$$\text{Score}(\emptyset, \text{▲▲▲}) + \text{Ins5}(\bullet, \circ)$$

si \bullet une base et \circ une paire

$$\text{Score}(\emptyset, \text{▲▲▲}) + \text{Ins3}(\bullet, \circ)$$

si \bullet une base et \circ une paire

Un algorithme d'alignement (2/4)

$$\text{Score}(\underbrace{\text{▲▲▲}, \text{▲▲▲}, \dots, \text{▲▲▲}}_n) = \text{Max}$$

$$\text{Score}(\emptyset, \text{▲▲▲}, \dots, \text{▲▲▲}) + \text{DelB}(\bullet) \quad \text{si } \bullet \text{ est une base}$$

$$\text{Score}(\text{▲▲▲}, \text{▲▲▲}, \dots, \text{▲▲▲}) + \text{InsB}(\bullet) \quad \text{si } \bullet \text{ est une base}$$

$$\text{Score}(\emptyset, \text{▲▲▲}, \dots, \text{▲▲▲}) + \text{SubB}(\bullet, \bullet) \quad \text{si } \bullet \text{ et } \bullet \text{ des bases}$$

$$\text{Score}(\text{▲▲▲}, \text{▲▲▲}, \dots, \text{▲▲▲}) + \text{DelP}(\bullet) \quad \text{si } \bullet \text{ est une paire}$$

$$\text{Score}(\text{▲▲▲}, \text{▲▲▲}, \dots, \text{▲▲▲}) + \text{Unpair}(\bullet, \bullet, \bullet) \quad \text{si } \bullet \text{ une paire, } \bullet \text{ et } \bullet \text{ des bases}$$

$$\text{Score}(\text{▲▲▲}, \text{▲▲▲}, \dots, \text{▲▲▲}) + \text{Del5}(\bullet, \bullet) \quad \text{si } \bullet \text{ une paire et } \bullet \text{ une base}$$

$$\text{Score}(\text{▲▲▲}, \text{▲▲▲}, \dots, \text{▲▲▲}) + \text{Del3}(\bullet, \bullet) \quad \text{si } \bullet \text{ une paire et } \bullet \text{ une base}$$

$$\text{Score}(\emptyset, \text{▲▲▲}, \dots, \text{▲▲▲}) + \text{Score}(\bullet, \text{▲▲▲}) \quad \text{si } \bullet \text{ une base et } \bullet \text{ une paire}$$

$$\text{Score}(\bullet, \text{▲▲▲}, \dots, \text{▲▲▲}) + \text{Score}(\emptyset, \text{▲▲▲}) \quad \text{si } \bullet \text{ une base et } \bullet \text{ une paire}$$

$$\text{Max}_{i+j=n} \left(\text{Score}(\emptyset, \text{▲▲▲}, \dots, \text{▲▲▲}) + \text{Score}(\text{▲▲▲}, \text{▲▲▲}, \dots, \text{▲▲▲}) \right) \quad \text{si } \bullet \text{ est une paire}$$

$$\text{Max}_{i+j=n} \left(\text{Score}(\text{▲▲▲}, \text{▲▲▲}, \dots, \text{▲▲▲}) + \text{Score}(\emptyset, \text{▲▲▲}, \dots, \text{▲▲▲}) \right) \quad \text{si } \bullet \text{ est une base}$$

Un algorithme d'alignement (3/4)

$$\text{Score}(\underbrace{\text{tree}_1 \dots \text{tree}_m}_{m}, \text{tree}_n) = \text{Max}$$

$$\text{Score}(\text{tree}_1 \dots \text{tree}_m, \text{tree}_n) + \text{DelB}(\bullet) \quad \text{si } \bullet \text{ est une base}$$

$$\text{Score}(\text{tree}_1 \dots \text{tree}_m, \emptyset) + \text{InsB}(\bullet) \quad \text{si } \bullet \text{ est une base}$$

$$\text{Score}(\text{tree}_1 \dots \text{tree}_m, \emptyset) + \text{SubB}(\bullet, \bullet) \quad \text{si } \bullet \text{ et } \bullet \text{ sont des bases}$$

$$\text{Score}(\text{tree}_1 \dots \text{tree}_m, \text{tree}_n) + \text{InsP}(\bullet) \quad \text{si } \bullet \text{ est une paire}$$

$$\text{Score}(\text{tree}_1 \dots \text{tree}_m, \text{tree}_n) + \text{Pair}(\bullet, \bullet, \bullet) \quad \text{si } \bullet \text{ et } \bullet \text{ des bases, et } \bullet \text{ une paire}$$

$$\text{Score}(\text{tree}_1 \dots \text{tree}_m, \text{tree}_n) + \text{Ins5}(\bullet, \bullet) \quad \text{si } \bullet \text{ une base et } \bullet \text{ une paire}$$

$$\text{Score}(\text{tree}_1 \dots \text{tree}_m, \text{tree}_n) + \text{Ins3}(\bullet, \bullet) \quad \text{si } \bullet \text{ une base et } \bullet \text{ une paire}$$

$$\text{Score}(\text{tree}_1 \dots \text{tree}_m, \emptyset) + \text{Score}(\text{tree}_n, \bullet) \quad \text{si } \bullet \text{ une paire et } \bullet \text{ une base}$$

$$\text{Score}(\text{tree}_1 \dots \text{tree}_m, \bullet) + \text{Score}(\text{tree}_n, \emptyset) \quad \text{si } \bullet \text{ une paire et } \bullet \text{ une base}$$

$$\text{Max}_{i+j=m} \left(\text{Score}(\text{tree}_1 \dots \text{tree}_i, \emptyset) + \text{Score}(\text{tree}_{i+1} \dots \text{tree}_m, \text{tree}_n) \right) \quad \left. \begin{array}{l} \text{si } \bullet \text{ est} \\ \text{une paire} \end{array} \right\}$$

$$\text{Max}_{i+j=m} \left(\text{Score}(\text{tree}_1 \dots \text{tree}_i, \text{tree}_n) + \text{Score}(\text{tree}_{i+1} \dots \text{tree}_m, \emptyset) \right) \quad \left. \begin{array}{l} \text{si } \bullet \text{ est} \\ \text{une base} \end{array} \right\}$$

Un algorithme d'alignement (4/4)

$$\text{Score}(\emptyset, \emptyset) = 0$$

$$\text{Score}(\begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \blacktriangle \quad \blacktriangle \end{array}, \emptyset) = \begin{cases} \text{Score}(\emptyset, \emptyset) + \text{DelB}(\bullet) & \text{si } \bullet \text{ est une base} \\ \text{Score}(\blacktriangle\blacktriangle\blacktriangle, \emptyset) + \text{DelP}(\bullet) & \text{si } \bullet \text{ est une paire} \end{cases}$$

$$\text{Score}(\blacktriangle \blacktriangle \begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \blacktriangle \quad \blacktriangle \end{array}, \emptyset) = \text{Score}(\blacktriangle \blacktriangle, \emptyset) + \text{Score}(\begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \blacktriangle \quad \blacktriangle \end{array}, \emptyset)$$

$$\text{Score}(\emptyset, \begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \blacktriangle \quad \blacktriangle \end{array}) = \begin{cases} \text{Score}(\emptyset, \emptyset) + \text{InsB}(\bullet) & \text{si } \bullet \text{ est une base} \\ \text{Score}(\emptyset, \blacktriangle\blacktriangle\blacktriangle) + \text{InsP}(\bullet) & \text{si } \bullet \text{ est une paire} \end{cases}$$

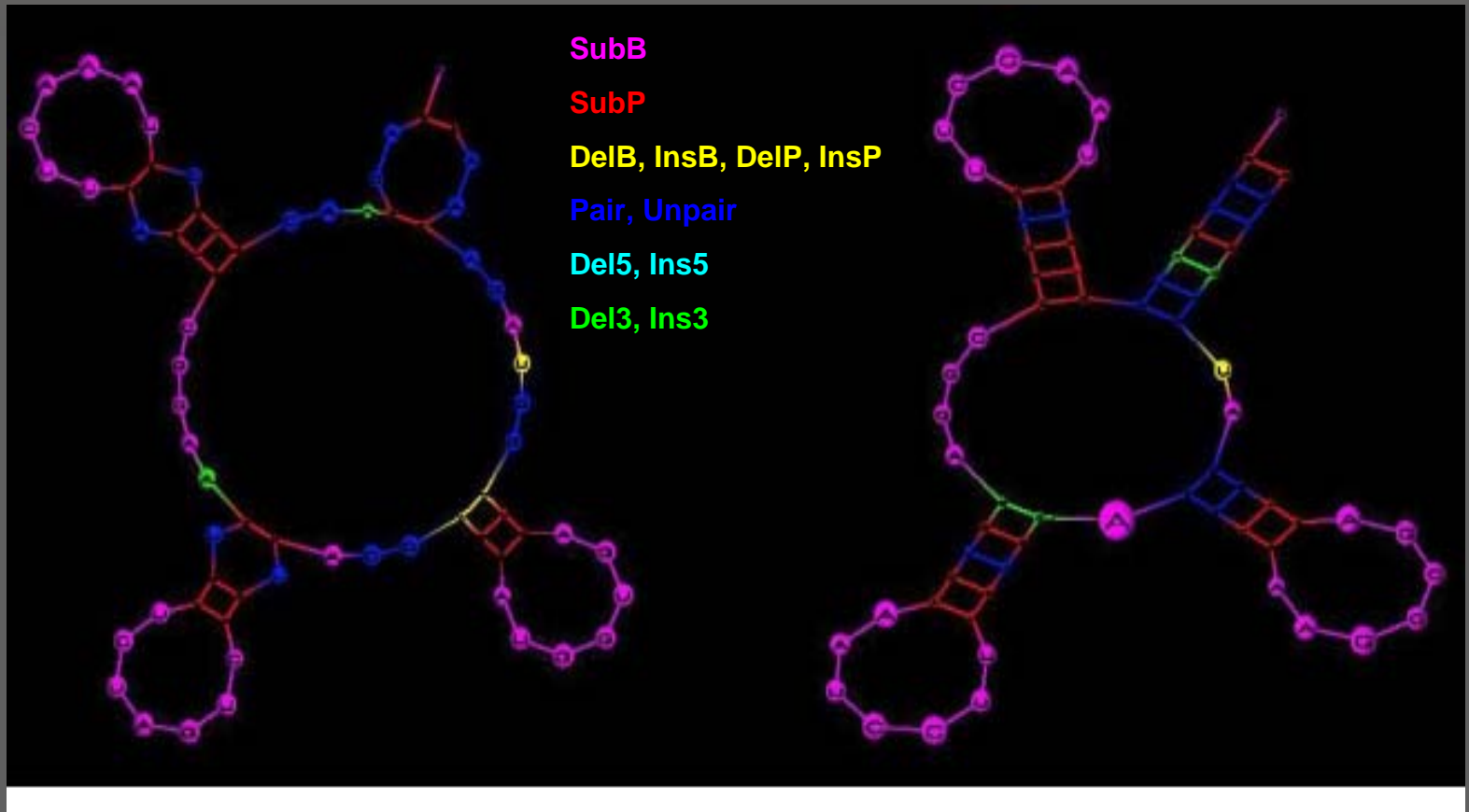
$$\text{Score}(\emptyset, \blacktriangle \blacktriangle \begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \blacktriangle \quad \blacktriangle \end{array}) = \text{Score}(\emptyset, \blacktriangle \blacktriangle) + \text{Score}(\emptyset, \begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \blacktriangle \quad \blacktriangle \end{array})$$

Edition et alignement d'arbres

	Schéma « arbres »	Schéma « ARN »
Edition	$O(n^3 \log n)$ [Zhang-Shasha 1989, Klein 1998]	NP-complet [Blin, Fertin, Sinoquet, Rusu 2003]
Alignement	$O(n^4)$ [Jiang, Wang, Zhang 1995]	$O(n^4)$ [Herrbach, AD, Dulucq, Touzet 2005]

Exemple : deux ARNt

Image avec Tulip (David Auber, LaBR)

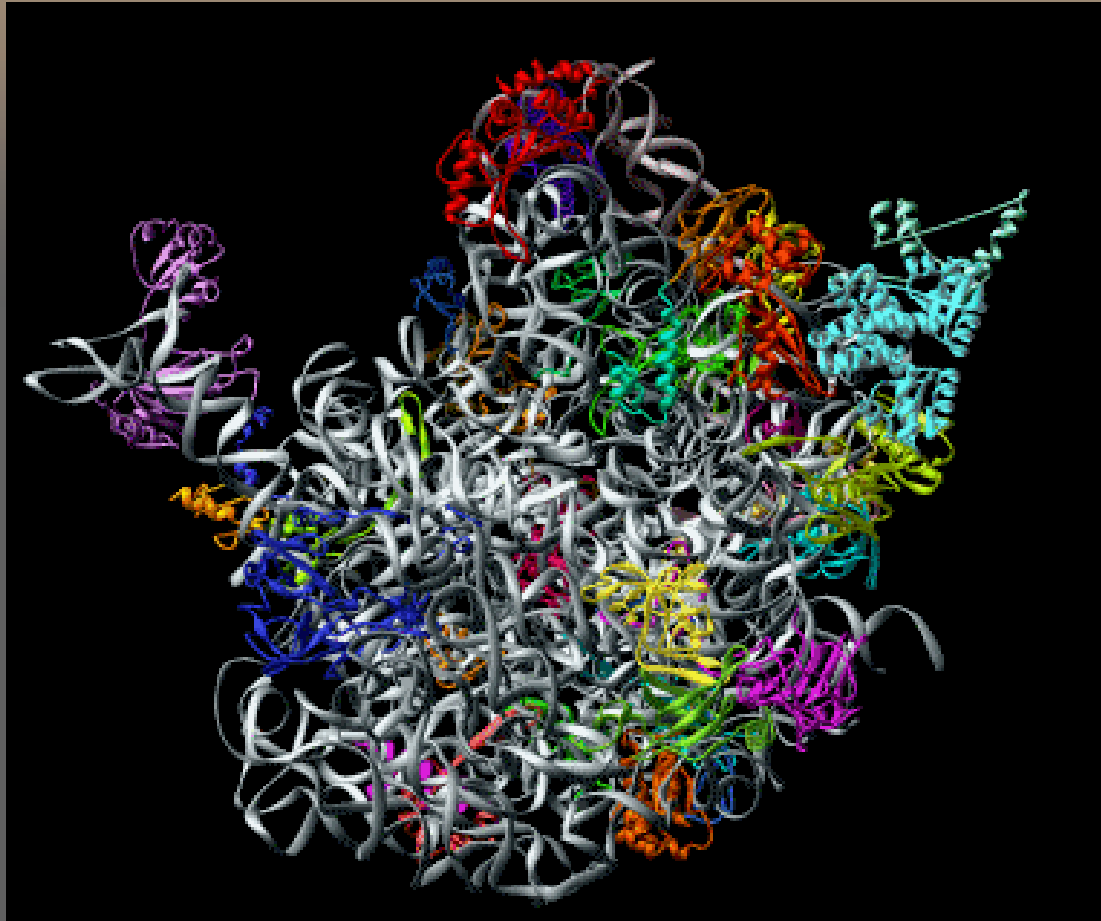


Homo sapiens

Bacillus subtilis

Crédits

- Serge Dulucq
- Claire Herrbach
- Yann Ponty
- Michel Termier
- Laurent Tichit
- Hélène Touzet
- Eric Westhof



navGraphe
ACI MdD