

## Cours M2 BIBS - Séance 2

### Équilibre de Boltzmann et comparaison

Yann Ponty

Bioinformatics Team  
École Polytechnique/CNRS/INRIA AMIB - France

13 Décembre 2010

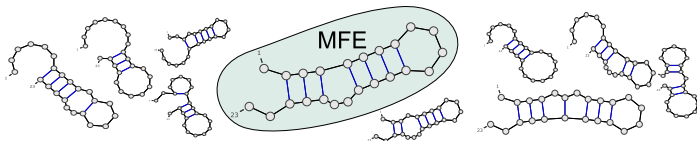
Yann Ponty Cours M2 BIBS - Séance 2 - Boltzmann et comparaison

## Ensemble canonique de Boltzmann

L'ARN *respire*  $\Rightarrow$  Il n'existe pas UNE unique conformation native.

Nouveau paradigme

Les conformations d'un ARN coexistent dans une distribution de Boltzmann.



**Conséquence :** La probabilité de la MFE peut être négligeable.  
 $\Rightarrow$  Comprendre les modes d'actions de l'ARN exige de prendre en considération l'ensemble des structures.

En particulier, des structures proches peuvent se *grouper* et devenir l'hypothèse la plus réaliste dans la recherche d'une conformation fonctionnelle.

Yann Ponty Cours M2 BIBS - Séance 2 - Boltzmann et comparaison

## Résumé

- 1 Ensemble de Boltzmann
  - Ensemble de Boltzmann
  - Nussinov : Minimisation  $\Rightarrow$  Comptage
  - Calcul de la fonction de partition
  - Échantillonnage statistique
- 2 Extensions
  - Validité d'un schéma
  - Structures sous-optimales
  - Pseudo-noeuds
- 3 Alignement et comparaison de structures d'ARN
  - Méthode géométrique
  - Alignement de structures secondaires
  - Méthodes hybrides

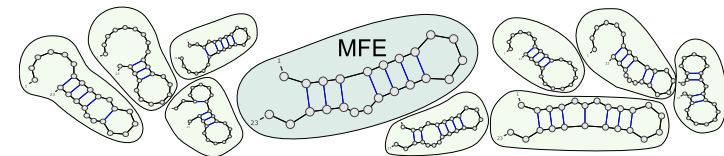
Yann Ponty Cours M2 BIBS - Séance 2 - Boltzmann et comparaison

## Ensemble canonique de Boltzmann

L'ARN *respire*  $\Rightarrow$  Il n'existe pas UNE unique conformation native.

Nouveau paradigme

Les conformations d'un ARN coexistent dans une distribution de Boltzmann.



**Conséquence :** La probabilité de la MFE peut être négligeable.  
 $\Rightarrow$  Comprendre les modes d'actions de l'ARN exige de prendre en considération l'ensemble des structures.

En particulier, des structures proches peuvent se *grouper* et devenir l'hypothèse la plus réaliste dans la recherche d'une conformation fonctionnelle.

Yann Ponty Cours M2 BIBS - Séance 2 - Boltzmann et comparaison

## Distribution de Boltzmann : Définition

Une distribution de Boltzmann pondère chaque structure  $S$  pour un ARN  $\omega$  par un **facteur de Boltzmann**  $\mathcal{B}_{S,\omega} = e^{\frac{-E_{S,\omega}}{RT}}$  où :

- $E_{S,\omega}$  est l'énergie libre de  $S$  (kCal.mol<sup>-1</sup>)
- $T$  est la température (K)
- $R$  est la constante des gaz parfaits (1.986.10<sup>-3</sup> kCal.K<sup>-1</sup>.mol<sup>-1</sup>)

Distribution renormalisée sur  $S_\omega$  par la **fonction de partition**

$$\mathcal{Z}_\omega = \sum_{S \in S_\omega} e^{\frac{-E_{S,\omega}}{RT}}$$

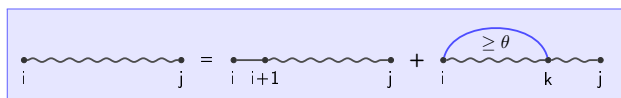
où  $S_\omega$  est l'ensemble des conformations compatibles avec  $\omega$ .

La **probabilité de Boltzmann** d'une structure  $S$  est alors donnée par

$$P_{S,\omega} = \frac{e^{\frac{-E_{S,\omega}}{RT}}}{\mathcal{Z}_\omega}$$

## Fonction de partition

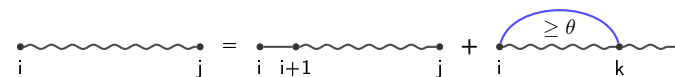
Fonction de partition = **Comptage pondéré** des structures compatibles



$$\mathcal{Z}_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

$$\mathcal{Z}_{i,j} = \sum \left\{ \begin{array}{l} \mathcal{Z}_{i+1,j} \\ \sum_{k=i+\theta+1}^j 1 \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \end{array} \right.$$

## Décomposition de Nussinov/Jacobson



Récurrance sur l'**énergie minimale** d'un repliement :

$$N_{i,t} = 0, \quad \forall t \in [i, i + \theta]$$

$$N_{i,j} = \min \left\{ \begin{array}{l} N_{i+1,j} \\ \min_{k=i+\theta+1}^j E_{i,k} + N_{i+1,k-1} + N_{k+1,j} \end{array} \right. \quad \begin{array}{l} (i \text{ non apparié}) \\ (i \text{ comp. avec } k) \end{array}$$

Récurrance de **comptage des structures compatibles** :

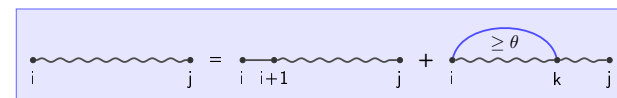
$$C_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

$$C_{i,j} = \sum \left\{ \begin{array}{l} C_{i+1,j} \\ \sum_{k=i+\theta+1}^j 1 \times C_{i+1,k-1} \times C_{k+1,j} \end{array} \right. \quad \begin{array}{l} (i \text{ non apparié}) \\ (i \text{ comp. avec } k) \end{array}$$

La décomposition est importante, le reste (MFE, comptage...) suit !

## Fonction de partition

Fonction de partition = **Comptage pondéré** des structures compatibles

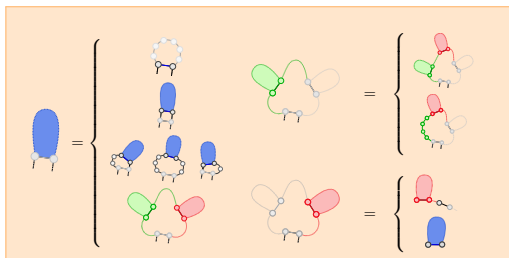


$$\mathcal{Z}_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

$$\mathcal{Z}_{i,j} = \sum \left\{ \begin{array}{l} \mathcal{Z}_{i+1,j} \\ \sum_{k=i+\theta+1}^j e^{\frac{-E_{bp}(i,k)}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \end{array} \right.$$

# Fonction de partition

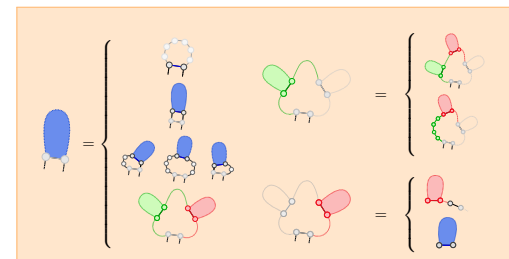
Fonction de partition = Comptage pondéré des structures compatibles



$$\begin{aligned} \mathcal{M}'_{i,j} &= \text{Min} \left\{ \begin{array}{l} E_H(i,j) \\ E_S(i,j) + \mathcal{M}'_{i+1,j-1} \\ \text{Min}(E_{BI}(i,i',j',j) + \mathcal{M}'_{i',j'}) \\ a + c + \text{Min}(\mathcal{M}_{i+1,k-1} + \mathcal{M}^1_{k,j-1}) \end{array} \right\} \\ \mathcal{M}_{i,j} &= \text{Min} \{ \text{Min}(\mathcal{M}_{i,k-1}, b(k-1)) + \mathcal{M}^1_{k,j} \} \\ \mathcal{M}^1_{i,j} &= \text{Min} \{ b + \mathcal{M}^1_{i,j-1}, c + \mathcal{M}'_{i,j} \} \end{aligned}$$

# Fonction de partition

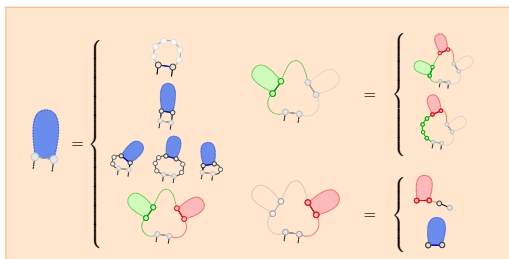
Fonction de partition = Comptage pondéré des structures compatibles



$$\begin{aligned} \mathcal{M}'_{i,j} &= \text{Min} \left\{ \begin{array}{l} e^{-\frac{E_H(i,j)}{RT}} \\ e^{-\frac{E_S(i,j)}{RT}} + \mathcal{M}'_{i+1,j-1} \\ \text{Min} \left( e^{-\frac{E_{BI}(i,i',j',j)}{RT}} + \mathcal{M}'_{i',j'} \right) \\ e^{-\frac{(a+c)}{RT}} + \text{Min}(\mathcal{M}_{i+1,k-1} + \mathcal{M}^1_{k,j-1}) \end{array} \right\} \\ \mathcal{M}_{i,j} &= \text{Min} \left\{ \text{Min}(\mathcal{M}_{i,k-1}, e^{-\frac{b(k-1)}{RT}}) + \mathcal{M}^1_{k,j} \right\} \\ \mathcal{M}^1_{i,j} &= \text{Min} \left\{ e^{-\frac{b}{RT}} + \mathcal{M}^1_{i,j-1}, e^{-\frac{c}{RT}} + \mathcal{M}'_{i,j} \right\} \end{aligned}$$

# Fonction de partition

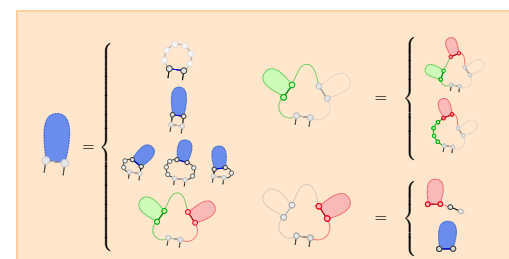
Fonction de partition = Comptage pondéré des structures compatibles



$$\begin{aligned} \mathcal{M}'_{i,j} &= \text{Min} \left\{ \begin{array}{l} e^{-\frac{E_H(i,j)}{RT}} \\ e^{-\frac{E_S(i,j)}{RT}} \mathcal{M}'_{i+1,j-1} \\ \text{Min} \left( e^{-\frac{E_{BI}(i,i',j',j)}{RT}} \mathcal{M}'_{i',j'} \right) \\ e^{-\frac{(a+c)}{RT}} \text{Min}(\mathcal{M}_{i+1,k-1}, \mathcal{M}^1_{k,j-1}) \end{array} \right\} \\ \mathcal{M}_{i,j} &= \text{Min} \left\{ \text{Min}(\mathcal{M}_{i,k-1}, e^{-\frac{b(k-1)}{RT}}) \mathcal{M}^1_{k,j} \right\} \\ \mathcal{M}^1_{i,j} &= \text{Min} \left\{ e^{-\frac{b}{RT}} \mathcal{M}^1_{i,j-1}, e^{-\frac{c}{RT}} \mathcal{M}'_{i,j} \right\} \end{aligned}$$

# Fonction de partition

Fonction de partition = Comptage pondéré des structures compatibles



$$\begin{aligned} \mathcal{Z}'(i,j) &= \sum \left\{ \begin{array}{l} e^{-\frac{E_H(i,j)}{RT}} \\ e^{-\frac{E_S(i,j)}{RT}} \mathcal{Z}'(i+1,j-1) \\ + \sum \left( e^{-\frac{E_{BI}(i,i',j',j)}{RT}} \mathcal{Z}'(i',j') \right) \\ + e^{-\frac{(a+c)}{RT}} \sum (\mathcal{Z}(i+1,k-1) \mathcal{Z}^1(k,j-1)) \end{array} \right\} \\ \mathcal{Z}(i,j) &= \sum (\mathcal{Z}(i,k-1) + e^{-\frac{b(k-1)}{RT}}) \mathcal{Z}^1(k,j) \\ \mathcal{Z}^1(i,j) &= e^{-\frac{b}{RT}} \mathcal{Z}^1(i,j-1) + e^{-\frac{c}{RT}} \mathcal{Z}'(i,j) \end{aligned}$$

## Fonction de partition

Fonction de partition = Comptage pondéré des structures compatibles

$$\mathcal{Z}_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

$$\mathcal{Z}_{i,j} = \sum \left\{ \sum_{k=i+\theta+1}^j \mathcal{Z}_{i+1,j} e^{-\frac{E_{bp}(i,k)}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \right.$$

Validité de la fonction de partition :

- Exhaustivité/non ambiguïté du schéma

## Fonction de partition

Fonction de partition = Comptage pondéré des structures compatibles

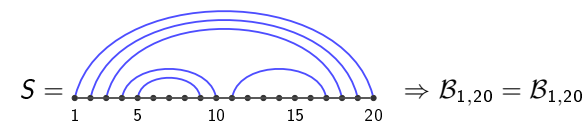
$$\mathcal{Z}_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

$$\mathcal{Z}_{i,j} = \sum \left\{ \sum_{k=i+\theta+1}^j \mathcal{Z}_{i+1,j} e^{-\frac{E_{bp}(i,k)}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \right.$$

Validité de la fonction de partition :

- Exhaustivité/non ambiguïté du schéma
- Correction du facteur de Boltzmann  
Facteur d'un backtrack = Produit des facteurs de ses parties  
Contributions énergétiques passent à l'exposant

Exemple :



## Fonction de partition

Fonction de partition = Comptage pondéré des structures compatibles

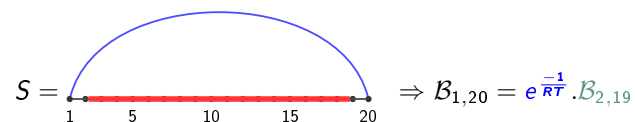
$$\mathcal{Z}_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

$$\mathcal{Z}_{i,j} = \sum \left\{ \sum_{k=i+\theta+1}^j \mathcal{Z}_{i+1,j} e^{-\frac{E_{bp}(i,k)}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \right.$$

Validité de la fonction de partition :

- Exhaustivité/non ambiguïté du schéma
- Correction du facteur de Boltzmann  
Facteur d'un backtrack = Produit des facteurs de ses parties  
Contributions énergétiques passent à l'exposant

Exemple :



## Fonction de partition

Fonction de partition = Comptage pondéré des structures compatibles

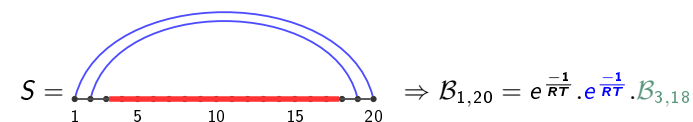
$$\mathcal{Z}_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

$$\mathcal{Z}_{i,j} = \sum \left\{ \sum_{k=i+\theta+1}^j \mathcal{Z}_{i+1,j} e^{-\frac{E_{bp}(i,k)}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \right.$$

Validité de la fonction de partition :

- Exhaustivité/non ambiguïté du schéma
- Correction du facteur de Boltzmann  
Facteur d'un backtrack = Produit des facteurs de ses parties  
Contributions énergétiques passent à l'exposant

Exemple :



## Fonction de partition

Fonction de partition = Comptage pondéré des structures compatibles

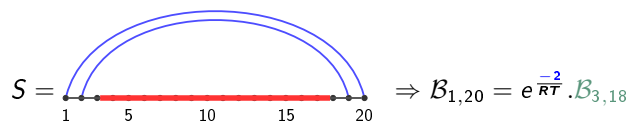
$$\mathcal{Z}_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

$$\mathcal{Z}_{i,j} = \sum \left\{ \sum_{k=i+\theta+1}^j e^{-\frac{E_{bp}(i,k)}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \right.$$

Validité de la fonction de partition :

- Exhaustivité/non ambiguïté du schéma
- Correction du facteur de Boltzmann  
Facteur d'un backtrack = Produit des facteurs de ses parties  
Contributions énergétiques passent à l'exposant

Exemple :



## Fonction de partition

Fonction de partition = Comptage pondéré des structures compatibles

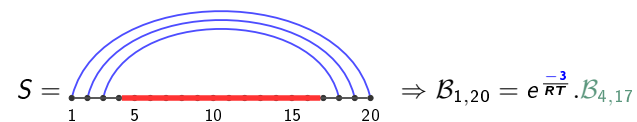
$$\mathcal{Z}_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

$$\mathcal{Z}_{i,j} = \sum \left\{ \sum_{k=i+\theta+1}^j e^{-\frac{E_{bp}(i,k)}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \right.$$

Validité de la fonction de partition :

- Exhaustivité/non ambiguïté du schéma
- Correction du facteur de Boltzmann  
Facteur d'un backtrack = Produit des facteurs de ses parties  
Contributions énergétiques passent à l'exposant

Exemple :



## Fonction de partition

Fonction de partition = Comptage pondéré des structures compatibles

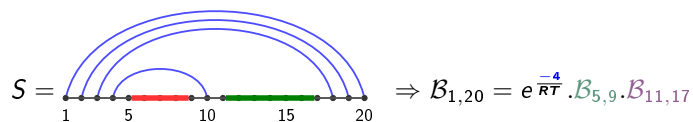
$$\mathcal{Z}_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

$$\mathcal{Z}_{i,j} = \sum \left\{ \sum_{k=i+\theta+1}^j e^{-\frac{E_{bp}(i,k)}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \right.$$

Validité de la fonction de partition :

- Exhaustivité/non ambiguïté du schéma
- Correction du facteur de Boltzmann  
Facteur d'un backtrack = Produit des facteurs de ses parties  
Contributions énergétiques passent à l'exposant

Exemple :



## Fonction de partition

Fonction de partition = Comptage pondéré des structures compatibles

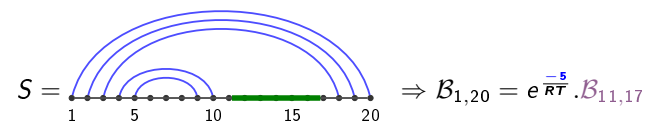
$$\mathcal{Z}_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

$$\mathcal{Z}_{i,j} = \sum \left\{ \sum_{k=i+\theta+1}^j e^{-\frac{E_{bp}(i,k)}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \right.$$

Validité de la fonction de partition :

- Exhaustivité/non ambiguïté du schéma
- Correction du facteur de Boltzmann  
Facteur d'un backtrack = Produit des facteurs de ses parties  
Contributions énergétiques passent à l'exposant

Exemple :



## Fonction de partition

Fonction de partition = Comptage pondéré des structures compatibles

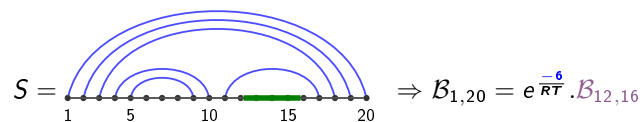
$$\mathcal{Z}_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

$$\mathcal{Z}_{i,j} = \sum \left\{ \sum_{k=i+\theta+1}^j e^{-\frac{E_{bp}(i,k)}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \right.$$

Validité de la fonction de partition :

- Exhaustivité/non ambiguïté du schéma
- Correction du facteur de Boltzmann  
Facteur d'un backtrack = Produit des facteurs de ses parties  
Contributions énergétiques passent à l'exposant

Exemple :



## Fonction de partition

Fonction de partition = Comptage pondéré des structures compatibles

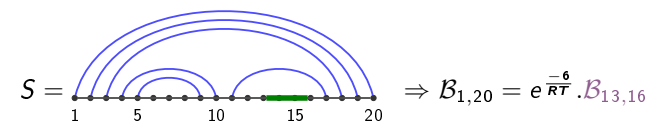
$$\mathcal{Z}_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

$$\mathcal{Z}_{i,j} = \sum \left\{ \sum_{k=i+\theta+1}^j e^{-\frac{E_{bp}(i,k)}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \right.$$

Validité de la fonction de partition :

- Exhaustivité/non ambiguïté du schéma
- Correction du facteur de Boltzmann  
Facteur d'un backtrack = Produit des facteurs de ses parties  
Contributions énergétiques passent à l'exposant

Exemple :



## Fonction de partition

Fonction de partition = Comptage pondéré des structures compatibles

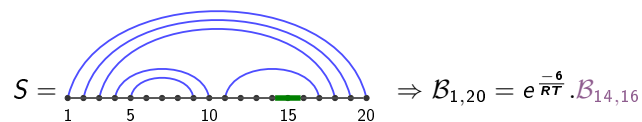
$$\mathcal{Z}_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

$$\mathcal{Z}_{i,j} = \sum \left\{ \sum_{k=i+\theta+1}^j e^{-\frac{E_{bp}(i,k)}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \right.$$

Validité de la fonction de partition :

- Exhaustivité/non ambiguïté du schéma
- Correction du facteur de Boltzmann  
Facteur d'un backtrack = Produit des facteurs de ses parties  
Contributions énergétiques passent à l'exposant

Exemple :



## Fonction de partition

Fonction de partition = Comptage pondéré des structures compatibles

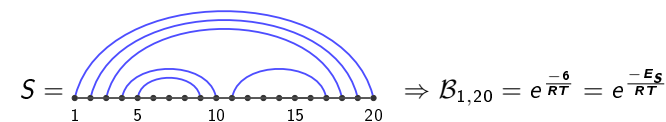
$$\mathcal{Z}_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

$$\mathcal{Z}_{i,j} = \sum \left\{ \sum_{k=i+\theta+1}^j e^{-\frac{E_{bp}(i,k)}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \right.$$

Validité de la fonction de partition :

- Exhaustivité/non ambiguïté du schéma
- Correction du facteur de Boltzmann  
Facteur d'un backtrack = Produit des facteurs de ses parties  
Contributions énergétiques passent à l'exposant

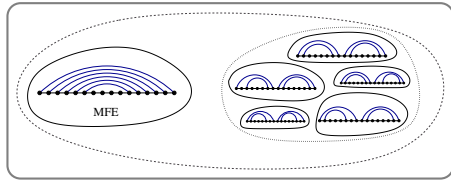
Exemple :



# Échantillonnage statistique de structures d'ARN

La MFE (Probabilité maximale) peut être *écrasée* par un ensemble  $\mathcal{B}$  de sous-optimaux structurellement similaires.

⇒ Conformation fonctionnelle trouvée plus probablement dans  $\mathcal{B}$ .



Expérience : [DCL05]

- Échantillonner des structures selon une probabilité de Boltzmann
- Effectuer un clustering
- Construire structure consensus dans le plus lourd cluster

⇒ Amélioration relative pour spécificité (+17.6%) et sensibilité (+21.74%, sauf Introns du groupe II)

## Problème

Comment engendrer des structures dans la distribution de Boltzmann ?

# Remontée stochastique

Algorithme (Reformulation SFold [DL03])

Précalcul : Calculer les matrices ( $\mathcal{Z}$ ,  $\mathcal{Z}'$ ,  $\mathcal{Z}^1$ ) des fonctions de partition.

Remontée stochastique :

- 1 Générer un nombre aléatoire  $r$  dans  $[0, \mathcal{Z}'(i, j))$

$$\mathcal{Z}'(i, j) = \sum \begin{cases} e^{-\frac{E_H(i, j)}{RT}} + e^{-\frac{E_S(i, j)}{RT}} \mathcal{Z}'(i+1, j-1) & \text{A} \\ \sum \left( e^{-\frac{E_{BJ}(i, i', j', j)}{RT}} \mathcal{Z}'(i', j') \right) & \text{B} \\ e^{-\frac{(a+c)}{RT}} \sum (\mathcal{Z}(i+1, k-1) \mathcal{Z}^1(k, j-1)) & \text{C} \end{cases}$$

# Remontée stochastique

Algorithme (Reformulation SFold [DL03])

Précalcul : Calculer les matrices ( $\mathcal{Z}$ ,  $\mathcal{Z}'$ ,  $\mathcal{Z}^1$ ) des fonctions de partition.

Remontée stochastique :

$$\mathcal{Z}'(i, j) = \begin{cases} \rightarrow e^{-\frac{E_H(i, j)}{RT}} + e^{-\frac{E_S(i, j)}{RT}} \mathcal{Z}'(i+1, j-1) & \text{A} \\ \rightarrow \sum \left( e^{-\frac{E_{BJ}(i, i', j', j)}{RT}} \mathcal{Z}'(i', j') \right) & \text{B} \\ \rightarrow e^{-\frac{(a+c)}{RT}} \sum (\mathcal{Z}(i+1, k-1) \mathcal{Z}^1(k, j-1)) & \text{C} \end{cases}$$

# Remontée stochastique

Algorithme (Reformulation SFold [DL03])

Précalcul : Calculer les matrices ( $\mathcal{Z}$ ,  $\mathcal{Z}'$ ,  $\mathcal{Z}^1$ ) des fonctions de partition.

Remontée stochastique :

- 1 Générer un nombre aléatoire  $r$  dans  $[0, \mathcal{Z}'(i, j))$

$$\mathcal{Z}'(i, j) = \sum \begin{cases} e^{-\frac{E_H(i, j)}{RT}} + e^{-\frac{E_S(i, j)}{RT}} \mathcal{Z}'(i+1, j-1) & \text{A} \\ \sum \left( e^{-\frac{E_{BJ}(i, i', j', j)}{RT}} \mathcal{Z}'(i', j') \right) & \text{B} \\ e^{-\frac{(a+c)}{RT}} \sum (\mathcal{Z}(i+1, k-1) \mathcal{Z}^1(k, j-1)) & \text{C} \end{cases}$$

$A_1 \mid A_2 \mid B_i \mid B_{i+1} \mid \dots \mid B_{j-1} \mid B_j \mid C_i \mid C_{i+1} \mid \dots \mid C_{j-1} \mid C_j$

# Remontée stochastique

Algorithme (Reformulation SFold [DL03])

Précalcul : Calculer les matrices ( $\mathcal{Z}$ ,  $\mathcal{Z}'$ ,  $\mathcal{Z}^1$ ) des fonctions de partition.

Remontée stochastique :

- 1 Générer un nombre aléatoire  $r$  dans  $[0, \mathcal{Z}'(i, j))$
- 2 Retirer à  $r$  les contributions à  $\mathcal{Z}'(i, j)$ , jusqu'à ce que  $r < 0$

$$\mathcal{Z}'(i, j) = \sum \left\{ \begin{array}{l} e^{-\frac{E_H(i, j)}{RT}} + e^{-\frac{E_S(i, j)}{RT}} \mathcal{Z}'(i+1, j-1) \quad \text{A} \\ \sum \left( e^{-\frac{E_{B_l}(i, i', j', j)}{RT}} \mathcal{Z}'(i', j') \right) \quad \text{B} \\ e^{-\frac{(\sigma+c)}{RT}} \sum (\mathcal{Z}(i+\frac{1}{r}, k-1) \mathcal{Z}^1(k, j-1)) \quad \text{C} \end{array} \right.$$

# Remontée stochastique

Algorithme (Reformulation SFold [DL03])

Précalcul : Calculer les matrices ( $\mathcal{Z}$ ,  $\mathcal{Z}'$ ,  $\mathcal{Z}^1$ ) des fonctions de partition.

Remontée stochastique :

- 1 Générer un nombre aléatoire  $r$  dans  $[0, \mathcal{Z}'(i, j))$
- 2 Retirer à  $r$  les contributions à  $\mathcal{Z}'(i, j)$ , jusqu'à ce que  $r < 0$

$$\mathcal{Z}'(i, j) = \sum \left\{ \begin{array}{l} e^{-\frac{E_H(i, j)}{RT}} + e^{-\frac{E_S(i, j)}{RT}} \mathcal{Z}'(i+1, j-1) \quad \text{A} \\ \sum \left( e^{-\frac{E_{B_l}(i, i', j', j)}{RT}} \mathcal{Z}'(i', j') \right) \quad \text{B} \\ e^{-\frac{(\sigma+c)}{RT}} \sum (\mathcal{Z}(i+\frac{1}{r}, k-1) \mathcal{Z}^1(k, j-1)) \quad \text{C} \end{array} \right.$$

# Remontée stochastique

Algorithme (Reformulation SFold [DL03])

Précalcul : Calculer les matrices ( $\mathcal{Z}$ ,  $\mathcal{Z}'$ ,  $\mathcal{Z}^1$ ) des fonctions de partition.

Remontée stochastique :

- 1 Générer un nombre aléatoire  $r$  dans  $[0, \mathcal{Z}'(i, j))$
- 2 Retirer à  $r$  les contributions à  $\mathcal{Z}'(i, j)$ , jusqu'à ce que  $r < 0$

$$\mathcal{Z}'(i, j) = \sum \left\{ \begin{array}{l} e^{-\frac{E_H(i, j)}{RT}} + e^{-\frac{E_S(i, j)}{RT}} \mathcal{Z}'(i+1, j-1) \quad \text{A} \\ \sum \left( e^{-\frac{E_{B_l}(i, i', j', j)}{RT}} \mathcal{Z}'(i', j') \right) \quad \text{B} \\ e^{-\frac{(\sigma+c)}{RT}} \sum (\mathcal{Z}(i+\frac{1}{r}, k-1) \mathcal{Z}^1(k, j-1)) \quad \text{C} \end{array} \right.$$

# Remontée stochastique

Algorithme (Reformulation SFold [DL03])

Précalcul : Calculer les matrices ( $\mathcal{Z}$ ,  $\mathcal{Z}'$ ,  $\mathcal{Z}^1$ ) des fonctions de partition.

Remontée stochastique :

- 1 Générer un nombre aléatoire  $r$  dans  $[0, \mathcal{Z}'(i, j))$
- 2 Retirer à  $r$  les contributions à  $\mathcal{Z}'(i, j)$ , jusqu'à ce que  $r < 0$

$$\mathcal{Z}'(i, j) = \sum \left\{ \begin{array}{l} e^{-\frac{E_H(i, j)}{RT}} + e^{-\frac{E_S(i, j)}{RT}} \mathcal{Z}'(i+1, j-1) \quad \text{A} \\ \sum \left( e^{-\frac{E_{B_l}(i, i', j', j)}{RT}} \mathcal{Z}'(i', j') \right) \quad \text{B} \\ e^{-\frac{(\sigma+c)}{RT}} \sum (\mathcal{Z}(i+\frac{1}{r}, k-1) \mathcal{Z}^1(k, j-1)) \quad \text{C} \end{array} \right.$$



## Remontée stochastique

Algorithme (Reformulation SFold [DL03])

Précalcul : Calculer les matrices ( $\mathcal{Z}$ ,  $\mathcal{Z}'$ ,  $\mathcal{Z}^1$ ) des fonctions de partition.

Remontée stochastique :

- 1 Générer un nombre aléatoire  $r$  dans  $[0, \mathcal{Z}'(i, j))$
- 2 Retirer à  $r$  les contributions à  $\mathcal{Z}'(i, j)$ , jusqu'à ce que  $r < 0$
- 3 Répéter sur les sous-structures

$$\mathcal{Z}'(i, j) = \sum \left\{ \begin{array}{l} e^{-\frac{E_H(i, j)}{RT}} + e^{-\frac{E_S(i, j)}{RT}} \mathcal{Z}'(i+1, j-1) \quad \text{A} \\ \sum \left( e^{-\frac{E_{B_l}(i, i', j', j)}{RT}} \mathcal{Z}'(i', j') \right) \quad \text{B} \\ e^{-\frac{(a+c)}{RT}} \sum (\mathcal{Z}(i+1, k-1) \mathcal{Z}^1(k, j-1)) \quad \text{C} \end{array} \right.$$

Correction : Chaque terme de la décomposition engendre

$\mathcal{T} \in \{\mathcal{A}_1, \dots, \mathcal{C}_j\}$ , et est choisi selon son facteur de Boltzmann cumulé  $\mathcal{B}(\mathcal{T})/\mathcal{Z} = \sum_{S \in \mathcal{T}} e^{-E/RT} / \mathcal{Z}$  (Par récurrence).

Chaque structure  $S \in \mathcal{S}_\omega$  est engendrée uniquement (Unambiguïté de Turner) par une séquence de choix d'ensembles

$\mathcal{S}_\omega \supset E_1 \supset E_2 \supset \dots \supset \{S\}$ .

La probabilité d'engendrer  $S$  est donc  $p_S = \frac{\mathcal{B}(E_1)}{\mathcal{B}(\mathcal{S}_\omega)} \cdot \frac{\mathcal{B}(E_2)}{\mathcal{B}(E_1)} \cdot \frac{\mathcal{B}(E_3)}{\mathcal{B}(E_2)} \dots \frac{\mathcal{B}(\{S\})}{\mathcal{B}(E_m)}$

## Remontée stochastique

Algorithme (Reformulation SFold [DL03])

Précalcul : Calculer les matrices ( $\mathcal{Z}$ ,  $\mathcal{Z}'$ ,  $\mathcal{Z}^1$ ) des fonctions de partition.

Remontée stochastique :

- 1 Générer un nombre aléatoire  $r$  dans  $[0, \mathcal{Z}'(i, j))$
- 2 Retirer à  $r$  les contributions à  $\mathcal{Z}'(i, j)$ , jusqu'à ce que  $r < 0$
- 3 Répéter sur les sous-structures

$$\mathcal{Z}'(i, j) = \sum \left\{ \begin{array}{l} e^{-\frac{E_H(i, j)}{RT}} + e^{-\frac{E_S(i, j)}{RT}} \mathcal{Z}'(i+1, j-1) \quad \text{A} \\ \sum \left( e^{-\frac{E_{B_l}(i, i', j', j)}{RT}} \mathcal{Z}'(i', j') \right) \quad \text{B} \\ e^{-\frac{(a+c)}{RT}} \sum (\mathcal{Z}(i+1, k-1) \mathcal{Z}^1(k, j-1)) \quad \text{C} \end{array} \right.$$

Correction : Chaque terme de la décomposition engendre

$\mathcal{T} \in \{\mathcal{A}_1, \dots, \mathcal{C}_j\}$ , et est choisi selon son facteur de Boltzmann cumulé  $\mathcal{B}(\mathcal{T})/\mathcal{Z} = \sum_{S \in \mathcal{T}} e^{-E/RT} / \mathcal{Z}$  (Par récurrence).

Chaque structure  $S \in \mathcal{S}_\omega$  est engendrée uniquement (Unambiguïté de Turner) par une séquence de choix d'ensembles

$\mathcal{S}_\omega \supset E_1 \supset E_2 \supset \dots \supset \{S\}$ .

La probabilité d'engendrer  $S$  est donc  $p_S = \frac{1}{\mathcal{B}(\mathcal{S}_\omega)} \cdot \frac{1}{1} \cdot \frac{1}{1} \dots \frac{\mathcal{B}(\{S\})}{1}$

## Remontée stochastique

Algorithme (Reformulation SFold [DL03])

Précalcul : Calculer les matrices ( $\mathcal{Z}$ ,  $\mathcal{Z}'$ ,  $\mathcal{Z}^1$ ) des fonctions de partition.

Remontée stochastique :

- 1 Générer un nombre aléatoire  $r$  dans  $[0, \mathcal{Z}'(i, j))$
- 2 Retirer à  $r$  les contributions à  $\mathcal{Z}'(i, j)$ , jusqu'à ce que  $r < 0$
- 3 Répéter sur les sous-structures

$$\mathcal{Z}'(i, j) = \sum \left\{ \begin{array}{l} e^{-\frac{E_H(i, j)}{RT}} + e^{-\frac{E_S(i, j)}{RT}} \mathcal{Z}'(i+1, j-1) \quad \text{A} \\ \sum \left( e^{-\frac{E_{B_l}(i, i', j', j)}{RT}} \mathcal{Z}'(i', j') \right) \quad \text{B} \\ e^{-\frac{(a+c)}{RT}} \sum (\mathcal{Z}(i+1, k-1) \mathcal{Z}^1(k, j-1)) \quad \text{C} \end{array} \right.$$

Correction : Chaque terme de la décomposition engendre

$\mathcal{T} \in \{\mathcal{A}_1, \dots, \mathcal{C}_j\}$ , et est choisi selon son facteur de Boltzmann cumulé  $\mathcal{B}(\mathcal{T})/\mathcal{Z} = \sum_{S \in \mathcal{T}} e^{-E/RT} / \mathcal{Z}$  (Par récurrence).

Chaque structure  $S \in \mathcal{S}_\omega$  est engendrée uniquement (Unambiguïté de Turner) par une séquence de choix d'ensembles

$\mathcal{S}_\omega \supset E_1 \supset E_2 \supset \dots \supset \{S\}$ .

La probabilité d'engendrer  $S$  est donc  $p_S = \frac{\mathcal{B}(\{S\})}{\mathcal{B}(\mathcal{S}_\omega)} = \frac{e^{-E_S/RT}}{\mathcal{Z}} = P_{S, \omega}$

## Complexité

Algorithme (Reformulation SFold [DL03])

- 1 Générer un nombre aléatoire  $r$  dans  $[0, \mathcal{Z}'(i, j))$
- 2 Retirer à  $r$  les contributions à  $\mathcal{Z}'(i, j)$ , jusqu'à ce que  $r < 0$
- 3 Répéter sur les sous-structures

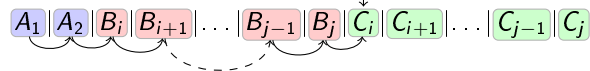
$$\mathcal{Z}'(i, j) = \sum \left\{ \begin{array}{l} \rightarrow e^{-\frac{E_H(i, j)}{RT}} + e^{-\frac{E_S(i, j)}{RT}} \mathcal{Z}'(i+1, j-1) \quad \text{A} \\ \rightarrow \sum \left( e^{-\frac{E_{B_l}(i, i', j', j)}{RT}} \mathcal{Z}'(i', j') \right) \quad \text{B} \\ \rightarrow e^{-\frac{(a+c)}{RT}} \sum (\mathcal{Z}(i+1, k-1) \mathcal{Z}^1(k, j-1)) \quad \text{C} \end{array} \right.$$

# Complexité

Algorithme (Reformulation SFold [DL03])

- 1 Générer un nombre aléatoire  $r$  dans  $[0, \mathcal{Z}'(i, j)]$
- 2 Retirer à  $r$  les contributions à  $\mathcal{Z}'(i, j)$ , jusqu'à ce que  $r < 0$
- 3 Répéter sur les sous-structures

$$\mathcal{Z}'(i, j) = \sum \left\{ \begin{array}{l} e^{-\frac{E_H(i, j)}{RT}} + e^{-\frac{E_S(i, j)}{RT}} \mathcal{Z}'(i+1, j-1) \quad \text{A} \\ \sum (e^{-\frac{E_{BJ}(i, i', j', j)}{RT}} \mathcal{Z}'(i', j')) \quad \text{B} \\ e^{-\frac{(a+c)}{RT}} \sum (\mathcal{Z}(i+1, k-1) \mathcal{Z}^1(k, j-1)) \quad \text{C} \end{array} \right.$$



Après  $\Theta(n)$  opérations, on réitère sur un interval de taille  $n-1$   
 $\Rightarrow$  Complexité du cas au pire en  $\mathcal{O}(n^2 k)$  pour  $k$  échantillons

Remarque : Instance pondérée d'un problème de génération aléatoire par la méthode *réursive* [Pon08].

Complexité en moyenne en  $\Theta(n\sqrt{n})$  dans l'hypothèse *tout appariement*.

Adaptation d'un parcours *Boustrophedon*  $\Rightarrow \mathcal{O}(n \log nk)$  au pire.

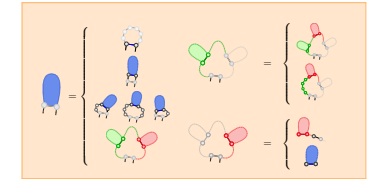
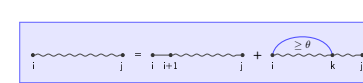
# Validité d'un schéma

Une preuve de correction possible :

Calcul correct localement

+ Toutes les conformations sont parcourues

$\Rightarrow$  Algorithme correct (Induction)



Forte certitude *mais* pas encore preuve (Séries génératrices).

# Validité d'un schéma

Une preuve de correction possible :

Calcul correct localement

+ Toutes les conformations sont parcourues

$\Rightarrow$  Algorithme correct (Induction)

$$C_{i,t} = 1, \forall t \in [i, i+\theta]$$

$$C_{i,j} = \sum \left\{ \begin{array}{l} C_{i+1,j} \\ \sum_{k=i+\theta+1}^j 1 \times C_{i+1,k-1} \times C_{k+1,j} \end{array} \right.$$

Homopolymère (Toute paire autorisée) +  $\theta = 1$   
 $\Rightarrow C_{1,n} = 1, 1, 1, 2, 4, 8, 17, 32, 82, 185, 423, \dots$



$$C'_{i,j} = \sum \left\{ \begin{array}{l} C'_{i+1,j-1} \\ \sum_{i',j'} C'_{i',j'} \\ \sum_k C_{i+1,k-1} \times C'_{k,j-1} \end{array} \right.$$

$$C_{i,j} = \sum_k ((C_{i,k-1} + 1) \times C'_{k,j})$$

$$C^1_{i,j} = C^1_{i,j-1} + C'_{i,j}$$

Homopolymère +  $\theta = 1$   
 $\Rightarrow C'_{1,n} = 0, 1, 1, 2, 4, 8, 17, 32, 82, 185, 423, \dots$

Forte certitude *mais* pas encore preuve (Séries génératrices).

# Repliement sous-optimal

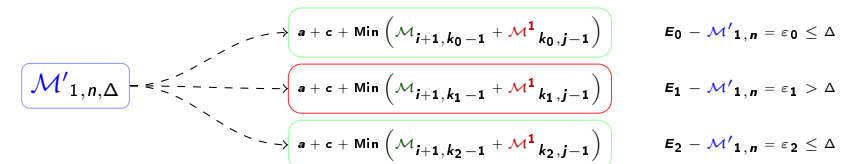
Prob. : Simplifications de l'énergie (Pseudo-noeuds, non-can.)

$\Rightarrow$  La structure *native* (fonctionnelle) pourrait être *ignorée*.

$\Rightarrow$  Engendrer des repliements sous-optimaux (RNASubopt [WFHS99]),

i.e. construire toutes les structures à  $\Delta$  KCal.mol<sup>-1</sup> de la MFE :

- Calculer la matrice des énergies minimales
- Effectuer un **Backtrack** sur toutes les contributions à  $\leq \Delta$  de la MFE



## Repliement sous-optimal

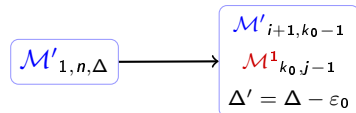
Prob. : Simplifications de l'énergie (Pseudo-nœuds, non-can.)

⇒ La structure *native* (fonctionnelle) pourrait être *ignorée*.

⇒ Engendrer des repliements sous-optimaux (RNASubopt [WFHS99]),

i.e. construire toutes les structures à  $\Delta$  KCal.mol<sup>-1</sup> de la MFE :

- Calculer la matrice des énergies minimales
- Effectuer un Backtrack sur toutes les contributions à  $\leq \Delta$  de la MFE
- Mettre à jour  $\Delta$  t.q. les futurs backtracks donnent  $\geq 1$  struct.



## Repliement sous-optimal

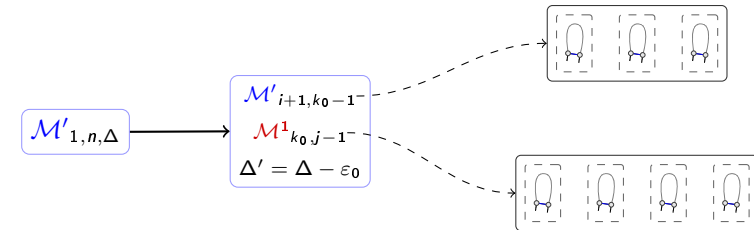
Prob. : Simplifications de l'énergie (Pseudo-nœuds, non-can.)

⇒ La structure *native* (fonctionnelle) pourrait être *ignorée*.

⇒ Engendrer des repliements sous-optimaux (RNASubopt [WFHS99]),

i.e. construire toutes les structures à  $\Delta$  KCal.mol<sup>-1</sup> de la MFE :

- Calculer la matrice des énergies minimales
- Effectuer un Backtrack sur toutes les contributions à  $\leq \Delta$  de la MFE
- Mettre à jour  $\Delta$  t.q. les futurs backtracks donnent  $\geq 1$  struct.
- Engendrer (Rec.) les sous-ensembles et combiner (brutal ou Tri)



## Repliement sous-optimal

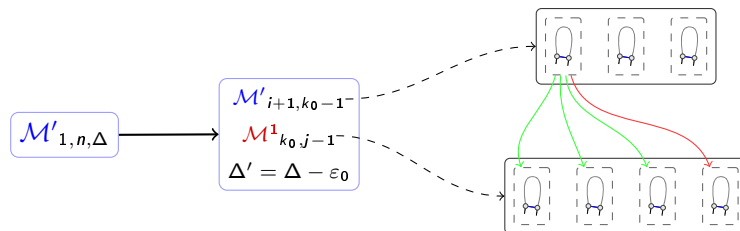
Prob. : Simplifications de l'énergie (Pseudo-nœuds, non-can.)

⇒ La structure *native* (fonctionnelle) pourrait être *ignorée*.

⇒ Engendrer des repliements sous-optimaux (RNASubopt [WFHS99]),

i.e. construire toutes les structures à  $\Delta$  KCal.mol<sup>-1</sup> de la MFE :

- Calculer la matrice des énergies minimales
- Effectuer un Backtrack sur toutes les contributions à  $\leq \Delta$  de la MFE
- Mettre à jour  $\Delta$  t.q. les futurs backtracks donnent  $\geq 1$  struct.
- Engendrer (Rec.) les sous-ensembles et combiner (brutal ou Tri)



## Repliement sous-optimal

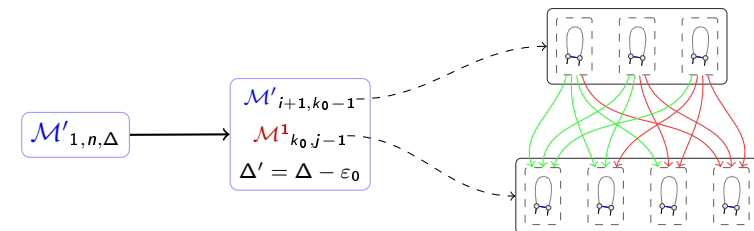
Prob. : Simplifications de l'énergie (Pseudo-nœuds, non-can.)

⇒ La structure *native* (fonctionnelle) pourrait être *ignorée*.

⇒ Engendrer des repliements sous-optimaux (RNASubopt [WFHS99]),

i.e. construire toutes les structures à  $\Delta$  KCal.mol<sup>-1</sup> de la MFE :

- Calculer la matrice des énergies minimales
- Effectuer un Backtrack sur toutes les contributions à  $\leq \Delta$  de la MFE
- Mettre à jour  $\Delta$  t.q. les futurs backtracks donnent  $\geq 1$  struct.
- Engendrer (Rec.) les sous-ensembles et combiner (brutal ou Tri)



## Repliement sous-optimal

Prob. : Simplifications de l'énergie (Pseudo-noeuds, non-can.)

⇒ La structure *native* (fonctionnelle) pourrait être *ignorée*.

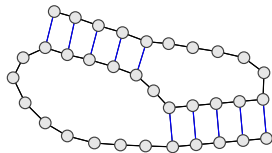
⇒ Engendrer des repliements sous-optimaux (RNASubopt [WFHS99]),  
i.e. construire toutes les structures à  $\Delta$  KCal.mol<sup>-1</sup> de la MFE :

- Calculer la matrice des énergies minimales
- Effectuer un Backtrack sur toutes les contributions à  $\leq \Delta$  de la MFE
- Mettre à jour  $\Delta$  t.q. les futurs backtracks donnent  $\geq 1$  struct.
- Engendrer (Rec.) les sous-ensembles et combiner (**brutal** ou **Tri**)

⇒ Complexité en temps (**Tri**) :  $\mathcal{O}(n^3 + nk \log(k))$   
( $k$  croît exponentiellement sur  $\Delta$ , mais bon...)

## Algorithme d'Akutsu/Uemura

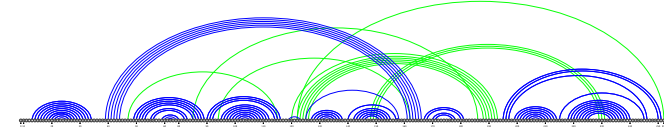
Le but est de capturer des catégories de pseudo-noeuds *simples*, mais très représentées.



**Idée** : Quand on retourne ce type de pseudonoeuds, il suffit de précalculer les meilleures configurations *en dessous* d'un triplet  $(i, j, k)$ , puis de regarder les configurations locales.

## Repliement avec pseudo-noeuds

Les pseudo-noeuds (et vrais noeuds) sont des constituants essentiels à la structuration de certains ARN.



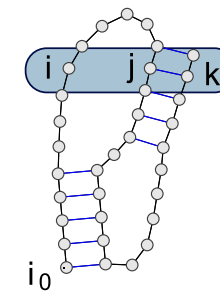
Ribozyme du groupe I

Leur absence historique au sein algorithmes de repliement est liée à la difficulté algorithmique des problèmes associés (Présence de pseudo-noeuds brise l'indépendance des repliements).

Type	Complexité	Référence
Structures secondaires	$\mathcal{O}(n^3)$	[MSZT99]
L&P	$\mathcal{O}(n^5)$	[LP00]
D&P	$\mathcal{O}(n^5)$	[DP03]
A&U	$\mathcal{O}(n^5)$	[Aku00]
R&E	$\mathcal{O}(n^6)$	[RE99]
Généraux	NP-complet	[LP00]

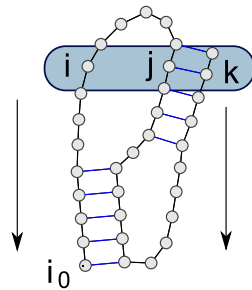
## Algorithme d'Akutsu/Uemura

Le but est de capturer des catégories de pseudo-noeuds *simples*, mais très représentées.



**Idée** : Quand on retourne ce type de pseudonoeuds, il suffit de précalculer les meilleures configurations *en dessous* d'un triplet  $(i, j, k)$ , puis de regarder les configurations locales.

Le but est de capturer des catégories de pseudo-nœuds *simples*, mais très représentée.



**Idée :** Quand on retourne ce type de pseudonœuds, il suffit de précalculer les meilleures configurations *en dessous* d'un triplet  $(i, j, k)$ , puis de regarder les configurations locales.

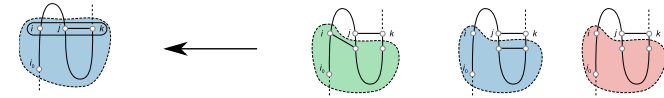
$r$  : Séquence d'ARN

$\Delta_{i,j}$  : Énergie de la paire  $(r_i, r_j)$ .



$$S_L^{i_0}(i, j, k) = \Delta_{i,j} + \min \left\{ \begin{array}{l} S_L^{i_0}(i-1, j+1, k), \\ S_M^{i_0}(i-1, j+1, k), \\ S_R^{i_0}(i-1, j+1, k) \end{array} \right\},$$

$$S_L^{i_0}(i, j, j) = \Delta_{i,j}, \quad \forall i < j$$

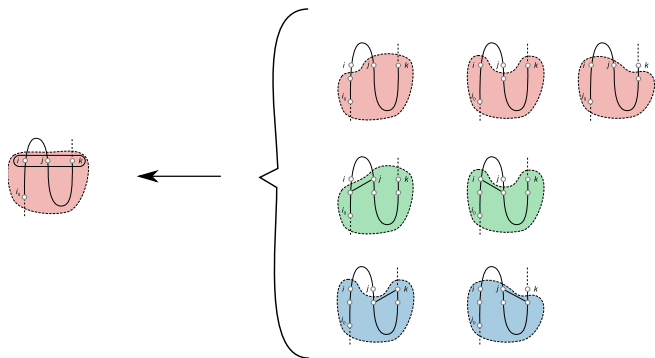


$$S_R^{i_0}(i, j, k) = \Delta_{j,k} + \min \left\{ \begin{array}{l} S_L^{i_0}(i, j+1, k-1), \\ S_M^{i_0}(i, j+1, k-1), \\ S_R^{i_0}(i, j+1, k-1) \end{array} \right\},$$

$$S_R^{i_0}(i_0-1, j, j+\theta+1) = \Delta_{j, j+\theta+1}, \quad \forall j$$

$r$  : Séquence d'ARN

$\Delta_{i,j}$  : Énergie de la paire  $(r_i, r_j)$ .



$$S_M^{i_0}(i, j, k) = \min \left\{ \begin{array}{l} S_M^{i_0}(i-1, j, k), S_M^{i_0}(i, j+1, k), S_M^{i_0}(i, j, k-1), \\ S_L^{i_0}(i-1, j, k), S_L^{i_0}(i, j+1, k) \\ S_R^{i_0}(i, j+1, k), S_R^{i_0}(i, j, k-1) \end{array} \right\}$$

$$S_L^{i_0}(i_0-1, j, k) = S_R^{i_0}(i_0-1, j, k) = S_M^{i_0}(i_0-1, j, k) = 0, \quad \forall j, k \text{ tq } k-j \leq \theta$$

$r$  : Séquence d'ARN

$\Delta_{i,j}$  : Énergie de la paire  $(r_i, r_j)$ .

L'équation générale sur le nombre de paires de bases pour la présence, sur l'intervalle  $(i_0, k_0)$ , de ce type de pseudo-nœuds est alors donné par

$$S_P(i_0, k_0) = \min_{i_0 \leq i < j < k \leq k_0} \left( S_L^{i_0}(i, j, k), S_R^{i_0}(i, j, k), S_M^{i_0}(i, j, k) \right)$$

On *insère* ces pseudonœuds au sein d'une structure secondaire classique au moyen d'une variante de Nussinov

$$S(i, j) = \min \left( S_P(i, j), S(i+1, j-1) + \Delta_{i,j}, \min_{i < k \leq j} (S(i, k-1) + S(k, j)) \right)$$

En utilisant une astuce dans l'ordre des calculs, on arrive à faire tomber la complexité à  $\mathcal{O}(n^4)$  dans un modèle de Nussinov, mais on reste en  $\mathcal{O}(n^5)$  dans le modèle de Turner.

- 1 Ensemble de Boltzmann
  - Ensemble de Boltzmann
  - Nussinov : Minimisation  $\Rightarrow$  Comptage
  - Calcul de la fonction de partition
  - Échantillonnage statistique
- 2 Extensions
  - Validité d'un schéma
  - Structures sous-optimales
  - Pseudo-noeuds
- 3 Alignement et comparaison de structures d'ARN
  - Méthode géométrique
  - Alignement de structures secondaires
  - Méthodes hybrides

## FR3D : Un exemple d'approche géométrique

Quand les structures tertiaires (3D) des ARN sont connues, le problème de l'alignement peut être abordé de façon **purement géométrique**.

### Problème

**Donnée :** Motif  $m$  et structure cible  $b$  (Ensembles de bases 3D).

**Résultat :** Matching de  $m$  et d'un sous-ensemble de  $b$  minimisant une **divergence** géométrique.

**Divergence géométrique :** Dans FR3D [SZS+08], une fonction  $D$  basée sur deux fonctions  $L$  et  $A$  d'erreur tenant compte respectivement de la superposabilité ( $L$ ) et de l'orientation des bases ( $A$ ) de  $m$  et  $b$ .

$$L = \sqrt{\min_{R,T} \sum_{i=1}^m \|b_i - R(T(m_i))\|^2} \quad A = \sqrt{\sum_{i=1}^m \alpha_i^2} \quad D = \frac{1}{m} \sqrt{L^2 + A^2}$$

$R, T$  : Rotation et translation.  $c_i$  : Barycentre pour la base  $m_i$ .  $\alpha_i$  : Écart entre les axes barycentre/bases dans  $m_i$  et  $b_i$ .

Exploration (Backtrack) + Élagage incrémental (Bornes sur  $D$ )  $\Rightarrow$  Explosion.

Mais recherche exacte pour des petits motifs.

## Pourquoi aligner structurellement des ARN

Une pression évolutive commune permet d'identifier une fonction commune.

Chez certains organismes (et pour certaines familles d'ARN), très faible conservation de la séquence. Cependant, la structure peut être bien plus conservée, et connue (Expérimentalement) ou déterminée par repliement.

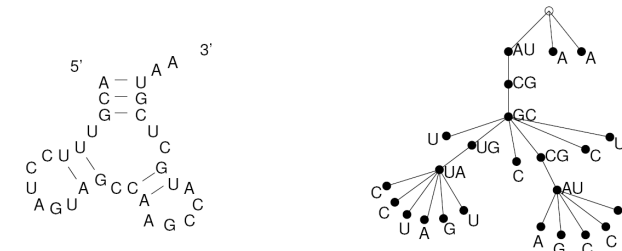
### Problèmes :

- **Édition :** Trouver la *distance* entre deux structures  $A$  et  $B$ .  
Quelle est la séquence d'opérations de coût minimal permettant de passer de  $A$  à  $B$  ? Déjà NP-complet pour deux structures secondaires [BFRS07].
- **Alignement :** Trouver une super-structure de coût minimal. Généralise la notion d'alignement de séquence. Polynomial pour des structures secondaires [BDD+08], NP-complet en 3D [SZS+08].  
**Variantes :** Alignement local ou global, Recherche de motifs.
- **Superposition :** Trouver une transformation géométrique (Rotation, translation, zoom) pour superposer *au mieux* les coordonnées de deux ARN de **matching connu**. Polynomial en 3D [McL82].

$\Rightarrow$  La difficulté algorithmique provient de la recherche du matching initial.

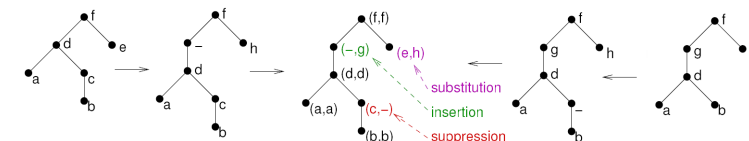
## Structures vers arbres

L'alignement de deux structures secondaires est basé sur une **représentation arborescente** de la structure secondaire<sup>1</sup>.



Paires de bases  $\Rightarrow$  noeuds internes      Bases non-appariées  $\Rightarrow$  Feuilles

**Alignement =** Construction d'un matching complet de coût minimal.



1. Illustrations empruntées à C. Herrbach

Alignement d'arbre<sup>2</sup>

$$\delta(\text{arbre}_1, \text{arbre}_2) = \min \begin{cases} \delta(\text{arbre}_1, \text{arbre}_2) + \text{del}(\bullet) \\ \delta(\text{arbre}_1, \text{arbre}_2) + \text{ins}(\bullet) \\ \delta(\text{arbre}_1, \text{arbre}_2) + \text{subst}(\bullet, \bullet) \end{cases}$$

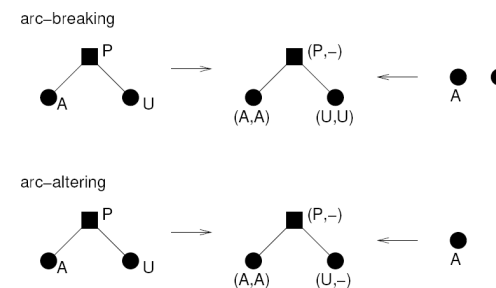
Alignement de forêt

$$\delta(\text{forêt}_1, \text{forêt}_2) = \min \begin{cases} \min\{\delta(\text{forêt}_1, \text{forêt}_2) + \delta(\text{forêt}_1, \text{forêt}_2) \mid \text{forêt}_1 = \text{forêt}_2\} + \text{del}(\bullet) \\ \min\{\delta(\text{forêt}_1, \text{forêt}_2) + \delta(\text{forêt}_1, \text{forêt}_2) \mid \text{forêt}_1 = \text{forêt}_2\} + \text{ins}(\bullet) \\ \delta(\text{forêt}_1, \text{forêt}_2) + \delta(\text{forêt}_1, \text{forêt}_2) \end{cases}$$

Complexité au pire en  $\mathcal{O}(n^4)$  [JWZ94], en moyenne en  $\mathcal{O}(n^2)$  [HDD07].  
Mais opérations spécifiques à l'ARN manquantes.

2. Idem

Basé sur l'algorithme de Jiang, Wang & Zhang + Intégrations d'opérations spécifiques à l'ARN<sup>3</sup>.



Possibilité de paramétrer les coûts des opérations, mais certaines opérations atomiques dans un modèle réaliste doivent être recomposées à partir des opérations disponibles. Par exemple, la substitution d'un sommet par une feuille est interdite directement.

3. Idem

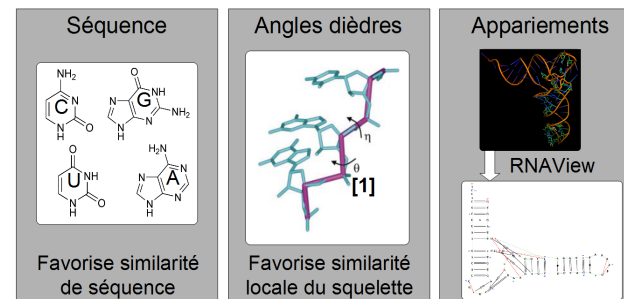
NestedAlign[BDD<sup>+</sup>08]

$$\delta(\text{arbre}_1, \text{arbre}_2) = \min \begin{cases} \delta(\text{arbre}_1, \text{arbre}_2) + \text{BDel}(\bullet) & \text{si } \bullet \text{ base} \\ \delta(\text{arbre}_1, \text{arbre}_2) + \text{BIns}(\bullet) & \text{si } \bullet \text{ base} \\ \delta(\text{arbre}_1, \text{arbre}_2) + \text{BSub}(\bullet, \bullet) & \text{si } \bullet \text{ et } \bullet \text{ bases} \\ \min\{\delta(\text{arbre}_1, \text{arbre}_2) + \delta(\text{arbre}_1, \text{arbre}_2) : \text{arbre}_1 = \text{arbre}_2\} + \text{PDel}(\bullet) & \text{si } \bullet \text{ paire} \\ \min\{\delta(\text{arbre}_1, \text{arbre}_2) + \delta(\text{arbre}_1, \text{arbre}_2) : \text{arbre}_1 = \text{arbre}_2\} + \text{PIns}(\bullet) & \text{si } \bullet \text{ paire} \\ \delta(\text{arbre}_1, \text{arbre}_2) + \delta(\text{arbre}_1, \text{arbre}_2) + \text{PSub}(\bullet, \bullet) & \text{si } \bullet \text{ et } \bullet \text{ paires} \\ \min\{\delta(\text{arbre}_1, \text{arbre}_2) + \delta(\text{arbre}_1, \text{arbre}_2) : \text{arbre}_1 = \text{arbre}_2\} + \text{Fus}(\bullet, \bullet, \bullet) & \text{si } \bullet \text{ paire et } \bullet \text{ base} \\ \min\{\delta(\text{arbre}_1, \text{arbre}_2) + \delta(\text{arbre}_1, \text{arbre}_2) : \text{arbre}_1 = \text{arbre}_2\} + \text{Sci}(\bullet, \bullet, \bullet) & \text{si } \bullet \text{ paire et } \bullet \text{ base} \\ \min\{\delta(\text{arbre}_1, \text{arbre}_2) + \delta(\text{arbre}_1, \text{arbre}_2) : \text{arbre}_1 = \text{arbre}_2\} + \text{GAlt}(\bullet, \bullet) & \text{si } \bullet \text{ paire et } \bullet \text{ base} \\ \min\{\delta(\text{arbre}_1, \text{arbre}_2) + \delta(\text{arbre}_1, \text{arbre}_2) : \text{arbre}_1 = \text{arbre}_2\} + \text{DAlt}(\bullet, \bullet) & \text{si } \bullet \text{ paire} \\ \min\{\delta(\text{arbre}_1, \text{arbre}_2) + \delta(\text{arbre}_1, \text{arbre}_2) : \text{arbre}_1 = \text{arbre}_2\} + \text{GComp}(\bullet, \bullet) & \text{si } \bullet \text{ paire et } \bullet \text{ base} \\ \min\{\delta(\text{arbre}_1, \text{arbre}_2) + \delta(\text{arbre}_1, \text{arbre}_2) : \text{arbre}_1 = \text{arbre}_2\} + \text{DComp}(\bullet, \bullet) & \text{si } \bullet \text{ paire} \end{cases}$$

Méthode hybride

DIAL [FPLC07] est une méthode hybride qui se concentre sur les comportements locaux.  
Idée : L'ARN est flexible, petite variation locale peuvent entraîner des grandes déviations géométriques.

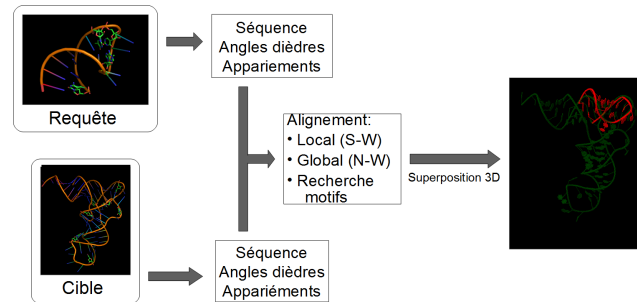
DIAL capture les similarités locales à trois niveau :



DIAL [FPLC07] est une méthode hybride qui se concentre sur les comportements locaux.

Idée : L'ARN est flexible, petite variation locale peuvent entraîner des grandes déviations géométriques.

Un algorithme d'alignement de séquence est alors utilisé



Tout dépend de ce que l'on a et veut :

- Modèle 3D :
  - Recherche d'un motif peu conservé en séquence : FR3D
  - Recherche d'un motif conservé : FR3D, DIAL ou DARTS
  - Recherche d'une structure entière : DIAL ou DARTS
- Structure secondaire :
  - Recherche d'un motif : NestedAlign
  - Alignement structure : RNAForester, NestedAlign

De nombreux autres programmes disponibles : Migal, Magnolia, ...  
 + Explosion des approches *par fragments* : FASTR3D, RNA FRABASE, ...

- Tatsuya Akutsu.  
**Dynamic programming algorithms for rna secondary structure prediction with pseudoknots.**  
*Discrete Appl. Math.*, 104(1-3) :45-62, 2000.
- G. Blin, A. Denise, S. Dulucq, C. Herrbach, and H. Touzet.  
**Alignment of rna structures.**  
*Transactions on Computational Biology and Bioinformatics*, ... .., 2008.  
 A paraître.
- Guillaume Blin, Guillaume Fertin, Irena Rusu, and Christine Sinoquet.  
**Extending the Hardness of RNA Secondary Structure Comparison.**  
 In Bo Chen, Mike Paterson, and Guochuan Zhang, editors, *ESCAPE'07*, volume 4614 of *LNCS*, pages 140-151, Hangzhou, China, Apr 2007.
- Y. Ding, C. Y. Chan, and C. E. Lawrence.  
**RNA secondary structure prediction by centroids in a boltzmann weighted ensemble.**  
*RNA*, 11 :1157-1166, 2005.
- Y. Ding and E. Lawrence.  
**A statistical sampling algorithm for RNA secondary structure prediction.**  
*Nucleic Acids Research*, 31(24) :7280-7301, 2003.
- Robert M Dirks and Niles A Pierce.  
**A partition function algorithm for nucleic acid secondary structure including pseudoknots.**  
*J Comput Chem*, 24(13) :1664-1677, Oct 2003.
- F. Ferrè, Y. Ponty, W. A. Lorenz, and Peter Clote.  
**Dial : A web server for the pairwise alignment of two RNA 3-dimensional structures using nucleotide, dihedral angle and base pairing similarities.**  
*Nucleic Acids Research*, 35(Web server issue) :W659-668, July 2007.

- Claire Herrbach, Alain Denise, and Serge Dulucq.  
**Average complexity of the jiang-wang-zhang pairwise tree alignment algorithm and of a rna secondary structure alignment algorithm.**  
 In *Proceedings of MACIS 2007, Second International Conference on Mathematical Aspects of Computer and Information Sciences*, 2007.
- M. Hochmann, B. Voss, and R. Giegerich.  
**Pure multiple RNA secondary structure alignments : A progressive profile approach.**  
*IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 01(1) :53-62, 2004.
- Tao Jiang, Lusheng Wang, and Kaizhong Zhang.  
**Alignment of trees - an alternative to tree edit.**  
 In *CPM '94 : Proceedings of the 5th Annual Symposium on Combinatorial Pattern Matching*, pages 75-86, London, UK, 1994. Springer-Verlag.
- R. B. Lyngsø and C. N. S. Pedersen.  
**RNA pseudoknot prediction in energy-based models.**  
*Journal of Computational Biology*, 7(3-4) :409-427, 2000.
- D. McLachlan.  
**Rapid comparison of protein structures.**  
*Acta crystallographica A*, 38(6) :871-873, 1982.
- D. H. Mathews, J. Sabina, M. Zuker, and D. H. Turner.  
**Expanded sequence dependence of thermodynamic parameters improves prediction of rna secondary structure.**  
*Journal of Molecular Biology*, 288(5) :911-940, May 1999.
- Y. Ponty.  
**Efficient sampling of RNA secondary structures from the boltzmann ensemble of low-energy : The boustrophedon method.**  
*Journal of Mathematical Biology*, 56(1-2) :107-127, Jan 2008.





E. Rivas and S. R. Eddy.

**A dynamic programming algorithm for rna structure prediction including pseudoknots.**  
*J Mol Biol*, 285(5) :2053–2068, Feb 1999.



M. Sarver, C. Zirbel, J. Stombaugh, A. Mokdad, and N. B. Leontis.

**FR3D : Finding local and composite recurrent structural motifs in RNA 3D.**  
*Journal of Mathematical Biology*, 56(1–2) :215–252, January 2008.



S. Wuchty, W. Fontana, I.L. Hofacker, and P. Schuster.

**Complete suboptimal folding of RNA and the stability of secondary structures.**  
*Biopolymers*, 49 :145–164, 1999.