

Cours M2 BIBS - Séance 2

Programmation dynamique avancée

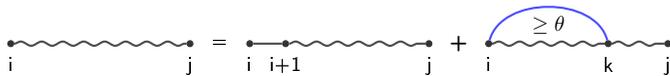
Yann Ponty

Bioinformatics Team
École Polytechnique/CNRS/INRIA AMIB – France

18 janvier 2010

Yann Ponty Cours M2 BIBS - Séance 1 - Repliement de l'ARN

Décomposition de Nussinov/Jacobson



Récurrance sur l'énergie minimale d'un repliement :

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

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

Récurrance de comptage des structures compatibles :

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

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

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

Yann Ponty Cours M2 BIBS - Séance 1 - Repliement de l'ARN

Résumé

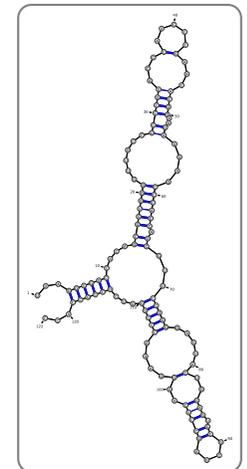
- 1 Rappels
 - Nussinov
- 2 Repliement - Modèle de Turner
 - Modèle de Turner
 - MFold/Unafold
 - Interlude : Validité d'un schéma
 - Structures sous-optimales
- 3 Changement de paradigme
 - Ensemble de Boltzmann
 - Calcul de la fonction de partition
 - Échantillonnage statistique

Yann Ponty Cours M2 BIBS - Séance 1 - Repliement de l'ARN

Modèle de Turner

Basée sur décomposition **non-ambiguë** en **boucles** de la structure 2^aire :

- Boucles internes
- Renflements
- Boucles terminales
- Boucles multiples
- Empilements



Énergies libres ΔG des boucles dépendent des bases, asymétrie, bases *libres* (dangle) ...

Déterminées expérimentalement
+ Interpolation pour grandes boucles.

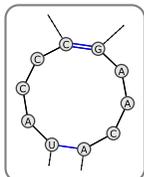
Meilleure résultats grâce à la prise en compte de l'empilement.

Yann Ponty Cours M2 BIBS - Séance 1 - Repliement de l'ARN

Modèle de Turner

Basée sur décomposition **non-ambiguë** en **boucles** de la structure 2^{aire} :

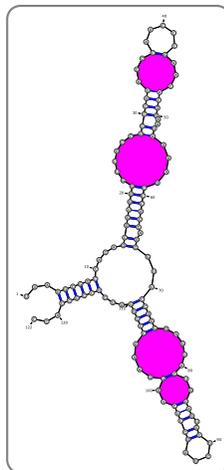
- Boucles internes
- Renflements
- Boucles terminales
- Boucles multiples
- Empilements



Énergies libres ΔG des boucles dépendent des bases, asymétrie, bases *libres* (dangle) ...

Déterminées expérimentalement
+ Interpolation pour grandes boucles.

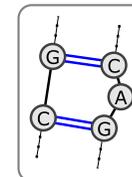
Meilleure résultats grâce à la prise en compte de l'empilement.



Modèle de Turner

Basée sur décomposition **non-ambiguë** en **boucles** de la structure 2^{aire} :

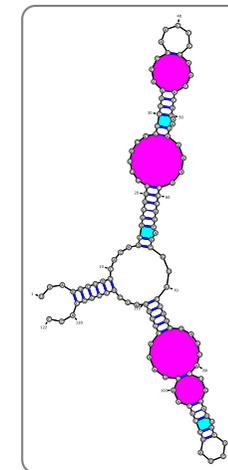
- Boucles internes
- Renflements
- Boucles terminales
- Boucles multiples
- Empilements



Énergies libres ΔG des boucles dépendent des bases, asymétrie, bases *libres* (dangle) ...

Déterminées expérimentalement
+ Interpolation pour grandes boucles.

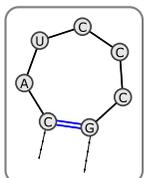
Meilleure résultats grâce à la prise en compte de l'empilement.



Modèle de Turner

Basée sur décomposition **non-ambiguë** en **boucles** de la structure 2^{aire} :

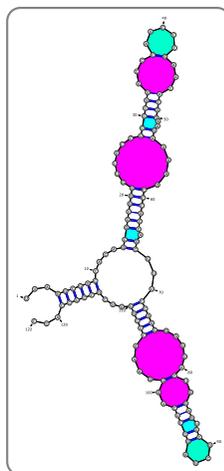
- Boucles internes
- Renflements
- Boucles terminales
- Boucles multiples
- Empilements



Énergies libres ΔG des boucles dépendent des bases, asymétrie, bases *libres* (dangle) ...

Déterminées expérimentalement
+ Interpolation pour grandes boucles.

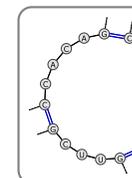
Meilleure résultats grâce à la prise en compte de l'empilement.



Modèle de Turner

Basée sur décomposition **non-ambiguë** en **boucles** de la structure 2^{aire} :

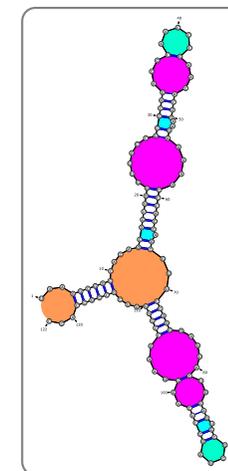
- Boucles internes
- Renflements
- Boucles terminales
- Boucles multiples
- Empilements



Énergies libres ΔG des boucles dépendent des bases, asymétrie, bases *libres* (dangle) ...

Déterminées expérimentalement
+ Interpolation pour grandes boucles.

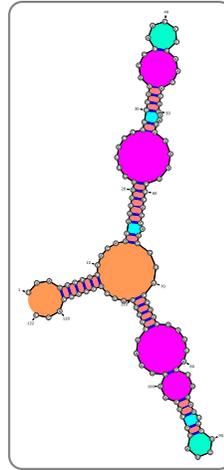
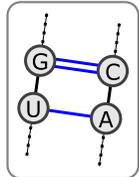
Meilleure résultats grâce à la prise en compte de l'empilement.



Modèle de Turner

Basée sur décomposition **non-ambiguë** en **boucles** de la structure 2^{aire} :

- Boucles internes
- Renflements
- Boucles terminales
- Boucles multiples
- Empilements

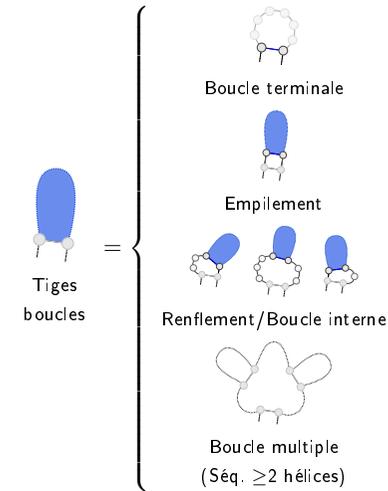


Énergies libres ΔG des boucles dépendent des bases, asymétrie, bases *libres* (dangle) ...

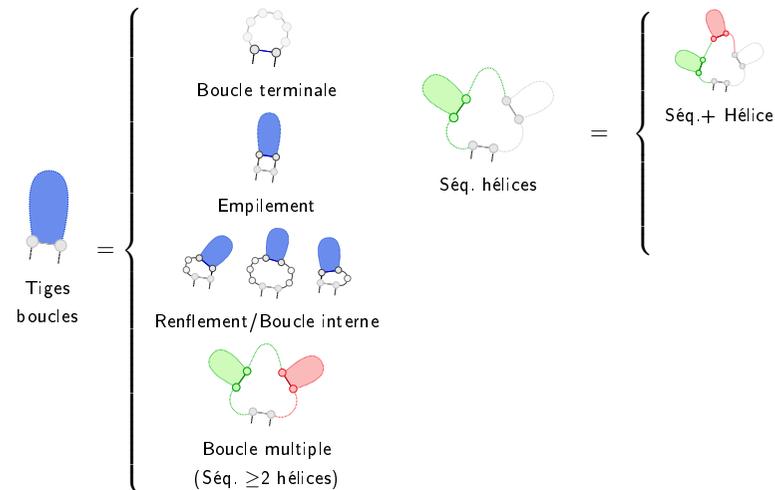
Déterminées expérimentalement
+ Interpolation pour grandes boucles.

Meilleure résultats grâce à la prise en compte de l'empilement.

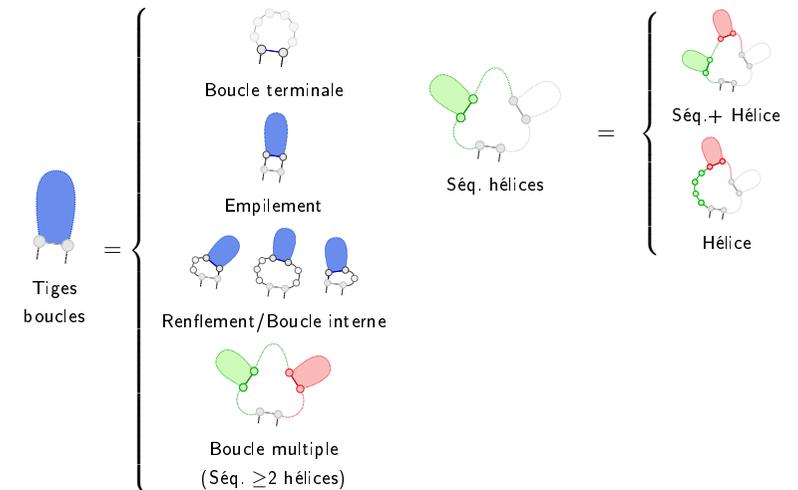
MFE DP equations



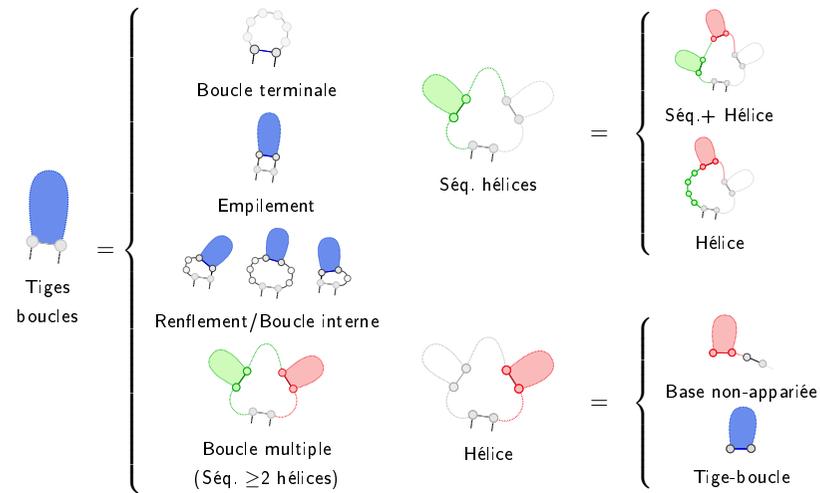
MFE DP equations



MFE DP equations

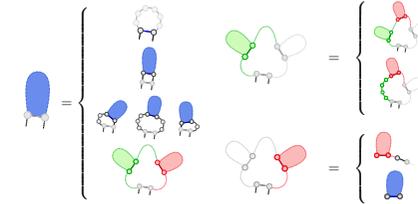


MFE DP equations



MFold Unafold

- $E_H(i, j)$: Energie de boucle terminale *fermée* par une paire (i, j)
- $E_{BI}(i, j)$: Energie de renflement ou boucle interne *fermée* par une paire (i, j)
- $E_S(i, j)$: Energie d'empilement $(i, j)/(i + 1, j - 1)$
- a, c, b : Pénalité de boucle multiple, hélice et non-appariées dans multiboucle.



Calcul des matrices

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

Remontée (Backtracking)

Reconstruction de la structure d'énergie minimale :

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

$\mathcal{O}(n)$ contributeurs potentiels au Min :

⇒ Complexité au pire en $\mathcal{O}(n^2)$ pour un backtrack naïf.

Garder les meilleures contributions aux Min ⇒ Backtrack en $\mathcal{O}(n)$

Complexités temps/mémoire en $\mathcal{O}(n^3)/\mathcal{O}(n^2)$ pour le précalcul¹

⇒ Unafold [MZ08] calcule la structure secondaire d'énergie minimale.

1. Avec une astuce pour les bulges/boucles internes ...

Remontée (Backtracking)

Reconstruction de la structure d'énergie minimale :

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

$\mathcal{O}(n)$ contributeurs potentiels au Min :

⇒ Complexité au pire en $\mathcal{O}(n^2)$ pour un backtrack naïf.

Garder les meilleures contributions aux Min ⇒ Backtrack en $\mathcal{O}(n)$

Complexités temps/mémoire en $\mathcal{O}(n^3)/\mathcal{O}(n^2)$ pour le précalcul¹

⇒ Unafold [MZ08] calcule la structure secondaire d'énergie minimale.

1. Avec une astuce pour les bulges/boucles internes ...

Remontée (Backtracking)

Reconstruction de la structure d'énergie minimale :

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

$$\mathcal{M}_{ij} = \text{Min}_k \{ \min (\mathcal{M}_{i, k-1}, b(k-1)) + \mathcal{M}^1_{k, j} \}$$

$$\mathcal{M}^1_{ij} = \text{Min}_k \{ b + \mathcal{M}^1_{i, j-1}, c + \mathcal{M}'_{ij} \}$$

$\mathcal{O}(n)$ contributeurs potentiels au Min :

⇒ Complexité au pire en $\mathcal{O}(n^2)$ pour un backtrack naïf.

Garder les meilleures contributions aux Min ⇒ Backtrack en $\mathcal{O}(n)$

Complexités temps/mémoire en $\mathcal{O}(n^3)/\mathcal{O}(n^2)$ pour le précalcul¹

⇒ `UnaFold` [MZ08] calcule la structure secondaire d'énergie minimale.

1. Avec une astuce pour les bulges/boucles internes ...

Remontée (Backtracking)

Reconstruction de la structure d'énergie minimale :

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

$$\mathcal{M}_{ij} = \text{Min}_k \{ \min (\mathcal{M}_{i, k-1}, b(k-1)) + \mathcal{M}^1_{k, j} \}$$

$$\mathcal{M}^1_{ij} = \text{Min}_k \{ b + \mathcal{M}^1_{i, j-1}, c + \mathcal{M}'_{ij} \}$$

$\mathcal{O}(n)$ contributeurs potentiels au Min :

⇒ Complexité au pire en $\mathcal{O}(n^2)$ pour un backtrack naïf.

Garder les meilleures contributions aux Min ⇒ Backtrack en $\mathcal{O}(n)$

Complexités temps/mémoire en $\mathcal{O}(n^3)/\mathcal{O}(n^2)$ pour le précalcul¹

⇒ `UnaFold` [MZ08] calcule la structure secondaire d'énergie minimale.

1. Avec une astuce pour les bulges/boucles internes ...

Remontée (Backtracking)

Reconstruction de la structure d'énergie minimale :

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

$$\mathcal{M}_{ij} \leftarrow \text{Min}_k \{ \min (\mathcal{M}_{i, k-1}, b(k-1)) + \mathcal{M}^1_{k, j} \}$$

$$\mathcal{M}^1_{ij} \leftarrow \text{Min}_k \{ b + \mathcal{M}^1_{i, j-1}, c + \mathcal{M}'_{ij} \}$$

$\mathcal{O}(n)$ contributeurs potentiels au Min :

⇒ Complexité au pire en $\mathcal{O}(n^2)$ pour un backtrack naïf.

Garder les meilleures contributions aux Min ⇒ Backtrack en $\mathcal{O}(n)$

Complexités temps/mémoire en $\mathcal{O}(n^3)/\mathcal{O}(n^2)$ pour le précalcul¹

⇒ `UnaFold` [MZ08] calcule la structure secondaire d'énergie minimale.

1. Avec une astuce pour les bulges/boucles internes ...

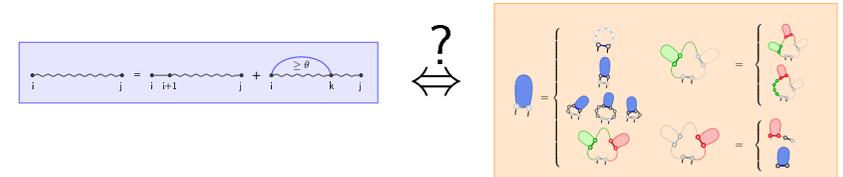
Validité d'un schéma

Une preuve de correction possible :

Calcul correct localement

+ Toutes les conformations sont parcourues

⇒ 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
- ⇒ Algorithme correct (Induction)

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

$$C_{i,j} = \sum_{k=i+\theta+1}^{C_{i+1,j}} 1 \times C_{i+1,k-1} \times C_{k+1,j}$$

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



$$C'_{i,j} = \sum \begin{cases} 1 \\ \sum_{i',j'} C'_{i',j'} \\ \sum_k C_{i+1,k-1} \times C_{k+1,j-1} \end{cases}$$

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

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

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

Forte certitude mais pas encore preuve (Séries génératrices).

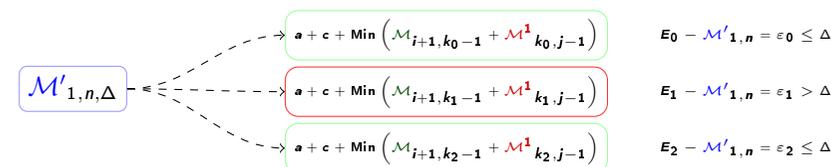
Repliement sous-optimal

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

⇒ La structure native (fonctionnelle) pourrait être ignorée.

⇒ **Engendrer des repliements sous-optimaux** (RNASubopt [WFHS99]),
 i.e. construire toutes les structures à Δ KCal.mol⁻¹ 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 Δ t.q. les futurs backtracks donnent ≥ 1 struct.
- Engendrer (Rec.) les sous-ensembles et combiner (**brutal** ou **Tri**)



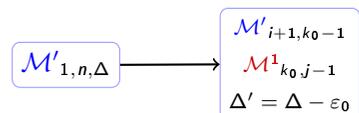
Repliement sous-optimal

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 Δ KCal.mol⁻¹ 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 Δ t.q. les futurs backtracks donnent ≥ 1 struct.**
- Engendrer (Rec.) les sous-ensembles et combiner (**brutal** ou **Tri**)



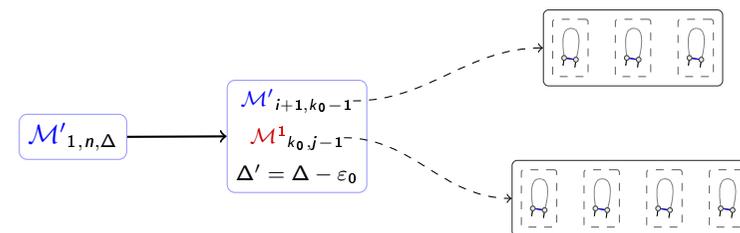
Repliement sous-optimal

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 Δ KCal.mol⁻¹ 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 Δ t.q. les futurs backtracks donnent ≥ 1 struct.
- **Engendrer (Rec.) les sous-ensembles et combiner (**brutal** ou **Tri**)**



Repliement sous-optimal

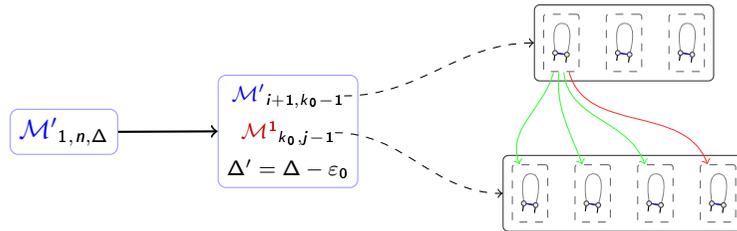
Prob. : Simplifications de l'énergie (Pseudo-noeuds, non-can.)

⇒ La structure native (fonctionnelle) pourrait être ignorée.

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

i.e. construire toutes les structures à Δ KCal.mol⁻¹ 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 Δ t.q. les futurs backtracks donnent ≥ 1 struct.
- Engendrer (Rec.) les sous-ensembles et combiner (brutal ou Tri)



Repliement sous-optimal

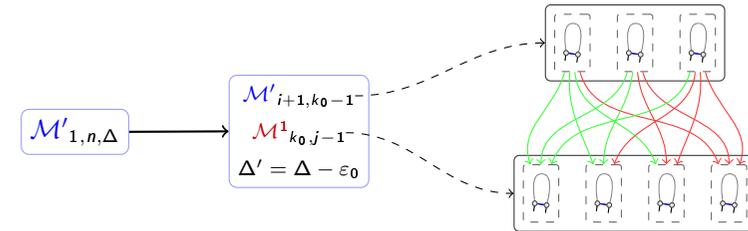
Prob. : Simplifications de l'énergie (Pseudo-noeuds, non-can.)

⇒ La structure native (fonctionnelle) pourrait être ignorée.

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

i.e. construire toutes les structures à Δ KCal.mol⁻¹ 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 Δ t.q. les futurs backtracks donnent ≥ 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 Δ KCal.mol⁻¹ 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 Δ t.q. les futurs backtracks donnent ≥ 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 Δ , mais bon...)

Résumé

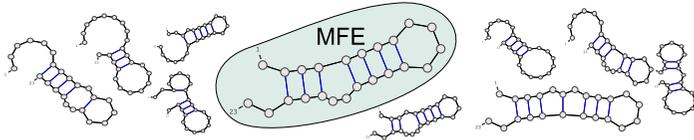
- 1 Rappels
 - Nussinov
- 2 Repliement - Modèle de Turner
 - Modèle de Turner
 - MFold/Unafold
 - Interlude : Validité d'un schéma
 - Structures sous-optimales
- 3 Changement de paradigme
 - Ensemble de Boltzmann
 - Calcul de la fonction de partition
 - Échantillonnage statistique

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.



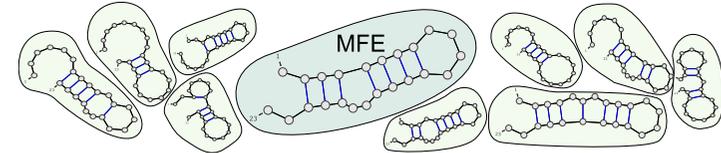
Conséquence : La probabilité de la MFE peut être négligeable.
 \Rightarrow La structure native (fonctionnelle) est alors à chercher dans les sous-optimales, où des structures très proches peuvent se *grouper*.

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 La structure native (fonctionnelle) est alors à chercher dans les sous-optimales, où des structures très proches peuvent se *grouper*.

Distribution de Boltzmann : Définitions

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

- $E_{S,\omega}$ est l'énergie libre de S (kCal.mol^{-1})
- T est la température (K)
- R est la constante des gaz parfaits ($1.986 \cdot 10^{-3} \text{ kCal.K}^{-1}.\text{mol}^{-1}$)

Distribution renormalisée sur S_ω par la fonction de partition

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

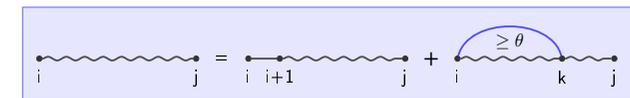
où S_ω est l'ensemble des conformations compatibles avec ω .

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

$$P_{S,\omega} = \frac{e^{-\frac{E_{S,\omega}}{RT}}}{Z_\omega}$$

Fonction de partition

Fonction de partition = Comptage pondéré des structures compatibles

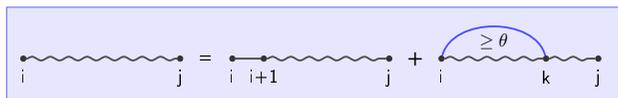


$$Z_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

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

Fonction de partition

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

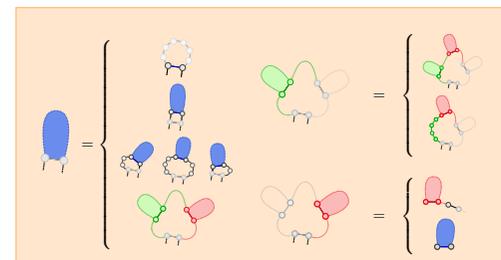


$$Z_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

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

Fonction de partition

Fonction de partition = Comptage pondéré des structures compatibles



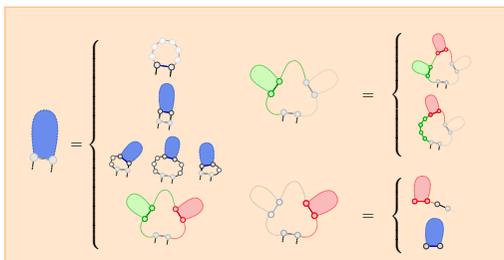
$$\mathcal{M}'_{i,j} = \text{Min} \left\{ \begin{array}{l} E_H(i,j) \\ E_S(i,j) + \mathcal{M}'_{i+1,j-1} \\ \text{Min}(E_{B_I}(i,i',j',j) + \mathcal{M}'_{i',j'}) \\ a + c + \text{Min}(\mathcal{M}_{i+1,k-1} + \mathcal{M}^1_{k,j-1}) \end{array} \right.$$

$$\mathcal{M}_{i,j} = \text{Min} \{ \text{Min}(\mathcal{M}_{i,k-1}, b(k-1)) + \mathcal{M}^1_{k,j} \}$$

$$\mathcal{M}^1_{i,j} = \text{Min} \{ b + \mathcal{M}^1_{i,j-1}, c + \mathcal{M}'_{i,j} \}$$

Fonction de partition

Fonction de partition = Comptage pondéré des structures compatibles



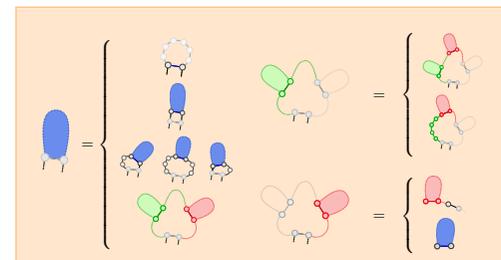
$$\mathcal{M}'_{i,j} = \text{Min} \left\{ \begin{array}{l} e^{-\frac{E_H(i,j)}{RT}} \\ e^{-\frac{E_S(i,j)}{RT}} + \mathcal{M}'_{i+1,j-1} \\ \text{Min} \left(e^{-\frac{E_{B_I}(i,i',j',j)}{RT}} + \mathcal{M}'_{i',j'} \right) \\ e^{-\frac{(a+c)}{RT}} + \text{Min}(\mathcal{M}_{i+1,k-1} + \mathcal{M}^1_{k,j-1})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

$$\mathcal{M}_{i,j} = \text{Min} \left\{ \text{Min}(\mathcal{M}_{i,k-1}, e^{-\frac{b(k-1)}{RT}}) + \mathcal{M}^1_{k,j} \right\}$$

$$\mathcal{M}^1_{i,j} = \text{Min} \left\{ e^{-\frac{b}{RT}} + \mathcal{M}^1_{i,j-1}, e^{-\frac{c}{RT}} + \mathcal{M}'_{i,j} \right\}$$

Fonction de partition

Fonction de partition = Comptage pondéré des structures compatibles



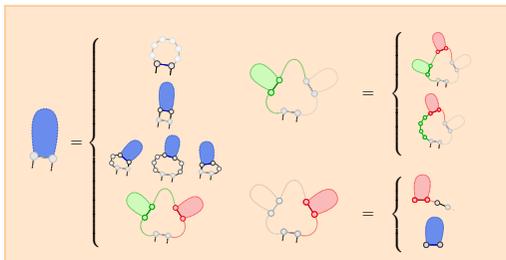
$$\mathcal{M}'_{i,j} = \text{Min} \left\{ \begin{array}{l} e^{-\frac{E_H(i,j)}{RT}} \\ e^{-\frac{E_S(i,j)}{RT}} + \mathcal{M}'_{i+1,j-1} \\ \text{Min} \left(e^{-\frac{E_{B_I}(i,i',j',j)}{RT}} + \mathcal{M}'_{i',j'} \right) \\ e^{-\frac{(a+c)}{RT}} \text{Min}(\mathcal{M}_{i+1,k-1}, \mathcal{M}^1_{k,j-1})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

$$\mathcal{M}_{i,j} = \text{Min} \left\{ \text{Min}(\mathcal{M}_{i,k-1}, e^{-\frac{b(k-1)}{RT}}) + \mathcal{M}^1_{k,j} \right\}$$

$$\mathcal{M}^1_{i,j} = \text{Min} \left\{ e^{-\frac{b}{RT}} + \mathcal{M}^1_{i,j-1}, e^{-\frac{c}{RT}} + \mathcal{M}'_{i,j} \right\}$$

Fonction de partition

Fonction de partition = Comptage pondéré des structures compatibles



$$\begin{aligned}
 Z'(i,j) &= \sum \left\{ \begin{aligned} &e^{-\frac{E_{\theta}(i,j)}{RT}} \\ &e^{-\frac{E_{\theta}(i,j)}{RT}} Z'(i+1, j-1) \\ &+ \sum \left(e^{-\frac{E_{\theta}(i',j')}{RT}} Z'(i', j') \right) \\ &+ e^{-\frac{E_{\theta}(i,j)}{RT}} \sum (Z(i+1, k-1) Z^1(k, j-1)) \end{aligned} \right. \\
 Z(i,j) &= \sum \left(Z(i, k-1) + e^{-\frac{E_{\theta}(i,j)}{RT}} \right) Z^1(k, j) \\
 Z^1(i,j) &= e^{-\frac{E_{\theta}(i,j)}{RT}} Z^1(i, j-1) + e^{-\frac{E_{\theta}(i,j)}{RT}} Z'(i, j)
 \end{aligned}$$

Fonction de partition

Fonction de partition = Comptage pondéré des structures compatibles

$$\begin{aligned}
 Z_{i,t} &= 1, \quad \forall t \in [i, i+\theta] \\
 Z_{i,j} &= \sum \left\{ \begin{aligned} &Z_{i+1,j} \\ &\sum_{k=i+\theta+1}^j e^{-\frac{E_{\theta}(i,j)}{RT}} \times Z_{i+1,k-1} \times Z_{k+1,j} \end{aligned} \right.
 \end{aligned}$$

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

Fonction de partition

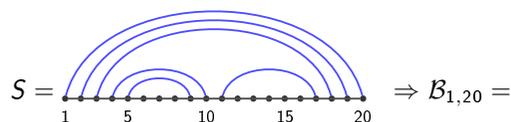
Fonction de partition = Comptage pondéré des structures compatibles

$$\begin{aligned}
 Z_{i,t} &= 1, \quad \forall t \in [i, i+\theta] \\
 Z_{i,j} &= \sum \left\{ \begin{aligned} &Z_{i+1,j} \\ &\sum_{k=i+\theta+1}^j e^{-\frac{E_{\theta}(i,j)}{RT}} \times Z_{i+1,k-1} \times Z_{k+1,j} \end{aligned} \right.
 \end{aligned}$$

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

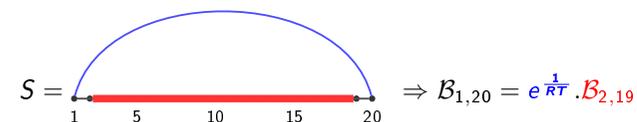
Fonction de partition = Comptage pondéré des structures compatibles

$$\begin{aligned}
 Z_{i,t} &= 1, \quad \forall t \in [i, i+\theta] \\
 Z_{i,j} &= \sum \left\{ \begin{aligned} &Z_{i+1,j} \\ &\sum_{k=i+\theta+1}^j e^{-\frac{E_{\theta}(i,j)}{RT}} \times Z_{i+1,k-1} \times Z_{k+1,j} \end{aligned} \right.
 \end{aligned}$$

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

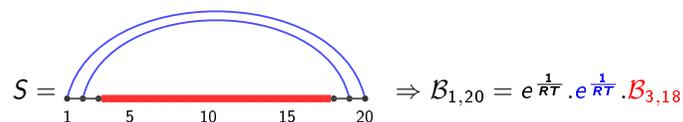
$$Z_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

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

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

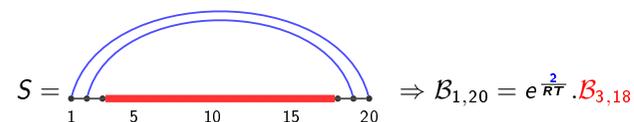
$$Z_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

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

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

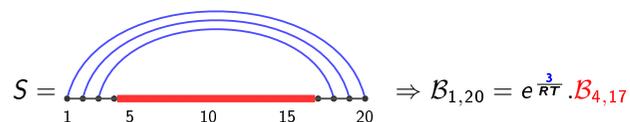
$$Z_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

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

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

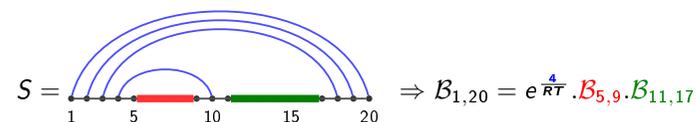
$$Z_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

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

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

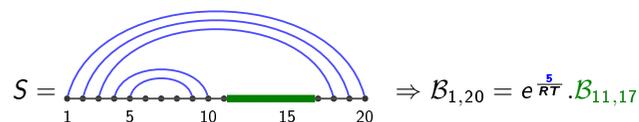
$$Z_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

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

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

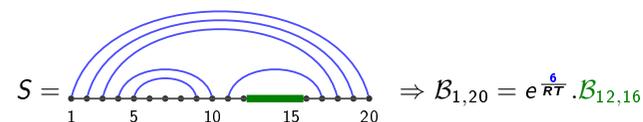
$$Z_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

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

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

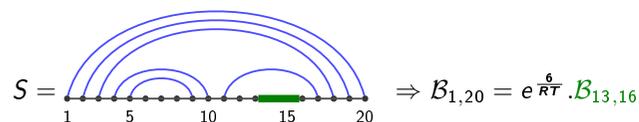
$$Z_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

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

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

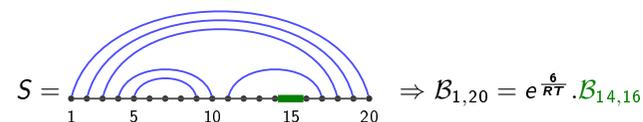
$$Z_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

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

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

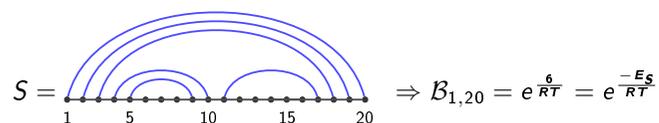
$$Z_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

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

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

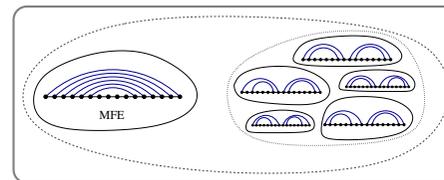
Exemple :



Échantillonnage statistique de structures d'ARN

Motivation de l'échantillonnage statistique

La MFE (Probabilité de Boltzmann max) \mathcal{M} peut être isolée et moins probable qu'un ensemble \mathcal{B} de sous-optimaux similaires structurellement. Alors, la structure native est plus probablement dans \mathcal{B} que \mathcal{M} [DCL05].



Approche :

- 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)

Remontée stochastique

Algorithme SFo1d [DL03] :

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

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

Remontée stochastique

Algorithme SFo1d [DL03] :

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

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

Remontée stochastique

Algorithme SFo1d [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

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

$\downarrow r$

Remontée stochastique

Algorithme SFo1d [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

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

$\downarrow r$

Remontée stochastique

Algorithme SFo1d [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

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

$\downarrow r$

Remontée stochastique

Algorithme SFo1d [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

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

$\downarrow r$

Algorithme SFo1d [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

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

Conclusion

Biologie/biochimie

- Structure aide à la compréhension de la fonction
 - Minimiser une fonction d'énergie n'est pas suffisant
 - Boltzmann = équilibre \Rightarrow Cinétique = convergence
- \Rightarrow Éviter une vision trop statique de la relation séquence/structure.

Bioinformatique

- Programmation dynamique (DP) est à la frontière du polynomial
 - DP résout exactement des problèmes complexes
 - Mais implémentation parfois problématiques
 - Séparer espace des conformations/énergie fait gagner du temps
- Chercher des ensembles de conformations discrets et simples².

2. Même si ça fait râler les biologistes/biochimistes ...

Algorithme SFo1d [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

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

Après $\Theta(n)$ opérations, on répète 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*. $\Rightarrow \mathcal{O}(n \log nk)$ au pire [Pon08].

Conclusion

Biologie/biochimie

- Structure aide à la compréhension de la fonction
 - Minimiser une fonction d'énergie n'est pas suffisant
 - Boltzmann = équilibre \Rightarrow Cinétique = convergence
- \Rightarrow Éviter une vision trop statique de la relation séquence/structure.

Bioinformatique

- Programmation dynamique (DP) est à la frontière du polynomial
 - DP résout exactement des problèmes complexes
 - Mais implémentation parfois problématiques
 - Séparer espace des conformations/énergie fait gagner du temps
- Chercher des ensembles de conformations discrets et simples².

2. Même si ça fait râler les biologistes/biochimistes ...

Conclusion

Biologie/biochimie

- Structure aide à la compréhension de la fonction
- Minimiser une fonction d'énergie n'est pas suffisant
- Boltzmann = équilibre \Rightarrow Cinétique = convergence

\Rightarrow Éviter une vision trop statique de la relation séquence/structure.

Bioinformatique

- Programmation dynamique (DP) est à la frontière du polynomial
- DP résout exactement des problèmes complexes
- Mais implémentation parfois problématiques
- Séparer espace des conformations/énergie fait gagner du temps

Chercher des ensembles de conformations discrets et simples².

2. Même si ça fait râler les biologistes/biochimistes ...

Conclusion

Biologie/biochimie

- Structure aide à la compréhension de la fonction
- Minimiser une fonction d'énergie n'est pas suffisant
- Boltzmann = équilibre \Rightarrow Cinétique = convergence

\Rightarrow Éviter une vision trop statique de la relation séquence/structure.

Bioinformatique

- Programmation dynamique (DP) est à la frontière du polynomial
- DP résout exactement des problèmes complexes
- Mais implémentation parfois problématiques
- Séparer espace des conformations/énergie fait gagner du temps

Chercher des ensembles de conformations discrets et simples².

2. Même si ça fait râler les biologistes/biochimistes ...

Conclusion

Biologie/biochimie

- Structure aide à la compréhension de la fonction
- Minimiser une fonction d'énergie n'est pas suffisant
- Boltzmann = équilibre \Rightarrow Cinétique = convergence

\Rightarrow Éviter une vision trop statique de la relation séquence/structure.

Bioinformatique

- Programmation dynamique (DP) est à la frontière du polynomial
- DP résout exactement des problèmes complexes
- Mais implémentation parfois problématiques
- Séparer espace des conformations/énergie fait gagner du temps

Chercher des ensembles de conformations discrets et simples².

2. Même si ça fait râler les biologistes/biochimistes ...

Conclusion

Biologie/biochimie

- Structure aide à la compréhension de la fonction
- Minimiser une fonction d'énergie n'est pas suffisant
- Boltzmann = équilibre \Rightarrow Cinétique = convergence

\Rightarrow Éviter une vision trop statique de la relation séquence/structure.

Bioinformatique

- Programmation dynamique (DP) est à la frontière du polynomial
- DP résout exactement des problèmes complexes
- Mais implémentation parfois problématiques
- Séparer espace des conformations/énergie fait gagner du temps

Chercher des ensembles de conformations discrets et simples².

2. Même si ça fait râler les biologistes/biochimistes ...

Conclusion

Biologie/biochimie

- Structure aide à la compréhension de la fonction
- Minimiser une fonction d'énergie n'est pas suffisant
- Boltzmann = équilibre \Rightarrow Cinétique = convergence

\Rightarrow Éviter une vision trop statique de la relation séquence/structure.

Bioinformatique

- Programmation dynamique (DP) est à la frontière du polynomial
- DP résout exactement des problèmes complexes
- Mais implémentation parfois problématiques
- Séparer espace des conformations/énergie fait gagner du temps

Chercher des ensembles de conformations discrets et simples².

2. Même si ça fait râler les biologistes/biochimistes ...

References |



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.



N. R. Markham and M. Zuker.
Bioinformatics, chapter UNAFold, pages 3–31.
Springer, 2008.



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.



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.