AMI2B Master - ABA Lecture

Boltzmann ensemble and partition function

Yann Ponty

Bioinformatics Team École Polytechnique/CNRS/INRIA AMIB – France

December 2nd, 2015

Outline

Focus on ab initio prediction

- Unambiguous Nussinov-like scheme
- Turner energy model
- Ab initio vs comparative

2 Boltzmann ensemble

- Nussinov : Minimisation ⇒ Comptage
- Calcul de la fonction de partition
- Échantillonnage statistique
- Inside/outside

Extensions

- Structures sous-optimales
- Pseudo-noeuds



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

$$N_{i,j} = \min \begin{cases} j & i \text{ unpaired} \\ \min_{k=i+\theta+1} \Delta G_{i,k} + N_{i+1,k-1} + N_{k+1,j} & i \text{ paired with } k \end{cases}$$



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

$$N_{i,j} = \min \begin{cases} j & i \text{ unpaired} \\ \min_{k=i+\theta+1} \Delta G_{i,k} + N_{i+1,k-1} + N_{k+1,j} & i \text{ paired with } k \end{cases}$$

Correctness. Goal = Show that MFE over interval [i, j] is indeed found in $N_{i,j}$ after completing the computation. Proceed by induction :

- ► Assume that property holds for any [i', j'] such that j' i' < n.
- ▶ Consider [i, j], j − i = n. Let MFE_{i,j} := Base-pairs of best struct. on [i, j]. Then first position i in MFE_{i,j} = is either :
 - ► Unpaired : MFE_{*i*,*j*} = MFE_{*i*+1,*j*} → free-energy = $N_{i+1,j}$
 - Paired to k : MFE_{1,j} = {(i, k)} ∪ MFE_{i+1,k-1} ∪ MFE_{k+1,j}. (Indeed, any BP between [i + 1, k − 1] and [k + 1, j] would cross (i, k))



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

$$N_{i,j} = \min \begin{cases} j & i \text{ unpaired} \\ \min_{k=i+\theta+1} \Delta G_{i,k} + N_{i+1,k-1} + N_{k+1,j} & i \text{ paired with } k \end{cases}$$

Correctness. Goal = Show that MFE over interval [i, j] is indeed found in $N_{i,j}$ after completing the computation. Proceed by induction :

- ► Assume that property holds for any [i', j'] such that j' i' < n.
- ► Consider [i, j], j − i = n. Let MFE_{i,j} := Base-pairs of best struct. on [i, j]. Then first position i in MFE_{i,j} = is either :
 - ► Unpaired : MFE_{*i*,*j*} = MFE_{*i*+1,*j*} \rightarrow free-energy = N_{*i*+1,*j*}
 - ▶ Paired to k: MFE_{*i*,*j*} = {(*i*, *k*)} \cup MFE_{*i*+1,*k*-1} \cup MFE_{*k*+1,*j*}. (Indeed, any BP between [*i* + 1, *k* - 1] and [*k* + 1, *j*] would cross (*i*, *k*)) \rightarrow free-energy = $\Delta G_{i,k} + N_{i+1,k-1} + N_{k+1,i}$



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

$$N_{i,j} = \min \begin{cases} j & i \text{ unpaired} \\ \min_{k=i+\theta+1} \Delta G_{i,k} + N_{i+1,k-1} + N_{k+1,j} & i \text{ paired with } k \end{cases}$$

Correctness. Goal = Show that MFE over interval [i, j] is indeed found in $N_{i,i}$ after completing the computation. Proceed by induction :

- Assume that property holds for any [i', j'] such that j' i' < n.
- Consider [i, j], j i = n. Let MFE_{*i*,*j*} := Base-pairs of best struct. on [i, j]. Then first position *i* in $MFE_{i,i} = is$ either :
 - ▶ Unpaired : $MFE_{i,j} = MFE_{i+1,j}$ \rightarrow fre ▶ Paired to k : $MFE_{i,j} = \{(i,k)\} \cup MFE_{i+1,k-1} \cup MFE_{k+1,j}$. \rightarrow free-energy = N_{i+1}
 - (Indeed, any BP between [i + 1, k 1] and [k + 1, i] would cross (i, k)) \rightarrow free-energy = $\Delta G_{i,k} + N_{i+1,k-1} + N_{k+1,i}$

Turner energy model

Based on unambiguous decomposition of 2ary structure into loops :

- Internal loops
- Bulges
- Terminal loops
- Multi loops
- Stackings

Free-energy Δ G of a loop depend on bases, assymmetry, dangles . . .

Experimentally determined + Interpolated for larger loops.



Turner energy model

Based on unambiguous decomposition of 2ary structure into loops :

Internal loops

- Bulges
- Terminal loops
- Multi loops
- Stackings



Free-energy Δ G of a loop depend on bases, assymmetry, dangles . . .

Experimentally determined + Interpolated for larger loops.



- Internal loops
- Bulges
- Terminal loops
- Multi loops
- Stackings



Free-energy Δ G of a loop depend on bases, assymmetry, dangles . . .

Experimentally determined + Interpolated for larger loops.



- Internal loops
- Bulges
- Terminal loops
- Multi loops
- Stackings



Free-energy Δ G of a loop depend on bases, assymmetry, dangles . . .

Experimentally determined + Interpolated for larger loops.



- Internal loops
- Bulges
- Terminal loops
- Multi loops
- Stackings



Free-energy Δ G of a loop depend on bases, assymmetry, dangles . . .

Experimentally determined + Interpolated for larger loops.



- Internal loops
- Bulges
- Terminal loops
- Multi loops
- Stackings



Free-energy Δ G of a loop depend on bases, assymmetry, dangles . . .

Experimentally determined + Interpolated for larger loops.











MFold Unafold

- E_H(i, j) : Energy of terminal loop enclosed by (i, j) pair
- E_{BI}(i, j) : Energy of bulge or internal loop enclosed by (i, j) pair
- $E_S(i,j)$: Energy of stacking (i,j)/(i+1,j-1)
- ▶ Penalty for multi loop (*a*), and occurrences of unpaired base (*b*) and helix (*c*) in multi loops.



DP recurrence

$$\mathcal{M}'_{i,j} = \min \begin{cases} E_{H}(i,j) \\ E_{S}(i,j) + \mathcal{M}'_{i+1,j-1} \\ \min_{j'} (E_{BJ}(i,j',j',j) + \mathcal{M}'_{i',j'}) \\ a + c + \min_{k} (\mathcal{M}_{i+1,k-1} + \mathcal{M}^{1}_{k,j-1}) \end{cases}$$

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

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

$$\mathcal{M}'_{i,j} = \operatorname{Min} \begin{cases} E_{\mathcal{H}}(i,j) \\ E_{S}(i,j) + \mathcal{M}'_{i+1,j-1} \\ Min_{i',j'}(E_{Bl}(i,i',j',j) + \mathcal{M}'_{i',j'}) \\ a + c + \operatorname{Min}_{k}(\mathcal{M}_{i+1,k-1} + \mathcal{M}^{1}_{k,j-1}) \\ \end{pmatrix} \\ \mathcal{M}_{i,j} = \operatorname{Min}_{k} \left\{ \min(\mathcal{M}_{i,k-1}, b(k-1)) + \mathcal{M}^{1}_{k,j} \right\} \\ \mathcal{M}^{1}_{i,j} = \operatorname{Min}_{k} \left\{ b + \mathcal{M}^{1}_{i,j-1}, c + \mathcal{M}'_{i,j} \right\} \end{cases}$$

Complexity :

For each min, $\mathcal{O}(n)$ potential contributors \Rightarrow **Worst-case** complexity in $\mathcal{O}(n^2)$ for **naive backtrack**. Keep best contributor for each Min \Rightarrow **Backtracking in** $\mathcal{O}(n)$

 \Rightarrow UnaFold [MZ08]/RNAFold [HFS⁺94] compute the MFE for the Turner model in **overall**¹ time/space complexities in $\mathcal{O}(n^3)/\mathcal{O}(n^2)$

$$\mathcal{M}'_{i,j} = \operatorname{Min}_{k,j} \left\{ \operatorname{Min}_{i,k-1}, b(k-1) \right\} + \mathcal{M}'_{k,j-1}$$

$$\mathcal{M}'_{i,j} = \operatorname{Min}_{k} \left\{ \operatorname{Min}_{i,k-1}, b(k-1) \right\} + \mathcal{M}'_{k,j-1}$$

Complexity :

For each min, $\mathcal{O}(n)$ potential contributors \Rightarrow **Worst-case** complexity in $\mathcal{O}(n^2)$ for **naive backtrack**. Keep best contributor for each Min \Rightarrow **Backtracking in** $\mathcal{O}(n)$

 \Rightarrow UnaFold [MZ08]/RNAFold [HFS⁺94] compute the MFE for the Turner model in **overall**¹ time/space complexities in $\mathcal{O}(n^3)/\mathcal{O}(n^2)$

$$\mathcal{M}'_{i,j} = \operatorname{Min}_{k} \left\{ b + \mathcal{M}^{1}_{i,j-1}, c + \mathcal{M}'_{i,j} \right\}$$

$$\mathcal{E}_{Bi}(i, j) + \mathcal{M}'_{i+1,j-1}$$

$$\mathcal{M}_{i,j'}(E_{Bi}(i, i', j', j) + \mathcal{M}'_{i',j'})$$

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

Complexity :

For each min, $\mathcal{O}(n)$ potential contributors \Rightarrow Worst-case complexity in $\mathcal{O}(n^2)$ for naive backtrack. Keep best contributor for each Min \Rightarrow Backtracking in $\mathcal{O}(n)$

⇒ UnaFold [MZ08]/RNAFold [HFS⁺94] compute the MFE for the Turner model in **overall**¹ time/space complexities in $\mathcal{O}(n^3)/\mathcal{O}(n^2)$

$$\mathcal{M}'_{i,j} = \operatorname{Min} \begin{cases} \mathbb{E}_{\mathcal{H}}(i,j) \\ \mathbb{E}_{S}(i,j) + \mathcal{M}'_{i+1,j-1} \\ \mathbb{M}_{i',j'}(\mathbb{E}_{Bl}(i,i',j',j) + \mathcal{M}'_{i',j'}) \\ \mathbb{A} + c + \operatorname{Min}_{k}(\mathcal{M}_{i+1,k-1} + \mathcal{M}^{1}_{k,j-1}) \\ \mathcal{M}_{i,j} = \operatorname{Min}_{k} \left\{ \min(\mathcal{M}_{i,k-1}, b(k-1)) + \mathcal{M}^{1}_{k,j} \right\} \\ \mathcal{M}^{1}_{i,j} = \operatorname{Min}_{k} \left\{ b + \mathcal{M}^{1}_{i,j-1}, c + \mathcal{M}'_{i,j} \right\}$$

Complexity :

For each min, $\mathcal{O}(n)$ potential contributors \Rightarrow Worst-case complexity in $\mathcal{O}(n^2)$ for naive backtrack. Keep best contributor for each Min \Rightarrow Backtracking in $\mathcal{O}(n)$

 \Rightarrow UnaFold [MZ08]/RNAFold [HFS⁺94] compute the MFE for the Turner model in **overall**¹ time/space complexities in $\mathcal{O}(n^3)/\mathcal{O}(n^2)$

$$\mathcal{M}'_{i,j} = \operatorname{Min} \begin{cases} \underbrace{E_{\mathcal{H}}(i,j)}_{E_{\mathcal{S}}(i,j) + \mathcal{M}'_{i+1,j-1}} \\ \underbrace{\operatorname{Min}_{i',j'}(E_{\mathcal{B}I}(i,i',j',j) + \mathcal{M}'_{i',j'})}_{a+c+\operatorname{Min}_{k}(\mathcal{M}_{i+1,k-1} + \mathcal{M}^{1}_{k,j-1})} \\ \underbrace{\mathcal{M}_{i,j}}_{mi,j} \leftarrow = -\operatorname{Min}_{k} \left\{ \min_{j=1}^{min} (\mathcal{M}_{i,k-1}, \mathcal{B}(\hat{k}-1)) + \mathcal{M}^{1}_{k,j} \right\} \\ \underbrace{\mathcal{M}^{1}_{i,j}}_{mi,j} \leftarrow = -\operatorname{Min}_{k} \left\{ -\mathcal{B} + \mathcal{M}^{1}_{i,j-1}; c + \widetilde{\mathcal{M}'}_{i,j} \right\} \end{cases}$$

Complexity :

For each min, $\mathcal{O}(n)$ potential contributors \Rightarrow Worst-case complexity in $\mathcal{O}(n^2)$ for naive backtrack. Keep best contributor for each Min \Rightarrow Backtracking in $\mathcal{O}(n)$

⇒ UnaFold [MZ08]/RNAFold [HFS⁺94] compute the MFE for the Turner model in **overall**¹ time/space complexities in $\mathcal{O}(n^3)/\mathcal{O}(n^2)$

$$\mathcal{M}'_{i,j} = \operatorname{Min} \begin{cases} E_{\mathcal{H}}(i,j) \\ E_{S}(i,j) + \mathcal{M}'_{i+1,j-1} \\ Min_{i',j'}(E_{Bl}(i,i',j',j) + \mathcal{M}'_{i',j'}) \\ a + c + \operatorname{Min}_{k}(\mathcal{M}_{i+1,k-1} + \mathcal{M}^{1}_{k,j-1}) \\ \mathcal{M}_{i,j} = \operatorname{Min}_{k} \left\{ \min(\mathcal{M}_{i,k-1}, b(k-1)) + \mathcal{M}^{1}_{k,j} \right\} \\ \mathcal{M}^{1}_{i,j} = \operatorname{Min}_{k} \left\{ b + \mathcal{M}^{1}_{i,j-1}, c + \mathcal{M}'_{i,j} \right\}$$

Complexity :

For each min, $\mathcal{O}(n)$ potential contributors \Rightarrow Worst-case complexity in $\mathcal{O}(n^2)$ for naive backtrack. Keep best contributor for each Min \Rightarrow Backtracking in $\mathcal{O}(n)$

 \Rightarrow UnaFold [MZ08]/RNAFold [HFS⁺94] compute the MFE for the Turner model in **overall**¹ time/space complexities in $\mathcal{O}(n^3)/\mathcal{O}(n^2)$

$$\mathcal{M}'_{i,j} = \operatorname{Min} \begin{cases} \underbrace{E_{\mathcal{H}}(i,j)}_{E_{\mathcal{S}}(i,j) + \mathcal{M}'_{i+1,j-1}} \\ \underbrace{\operatorname{Min}_{i',j'}(E_{\mathcal{B}I}(i,i',j',j) + \mathcal{M}'_{i',j'})}_{a+c+\operatorname{Min}_{k}(\mathcal{M}_{i+1,k-1} + \mathcal{M}^{1}_{k,j-1})} \\ \underbrace{\mathcal{M}_{i,j}}_{a,j} = \operatorname{Min}_{k} \left\{ \min(\mathcal{M}_{i,k-1}, b(k-1)) + \mathcal{M}^{1}_{k,j} \right\} \\ \underbrace{\mathcal{M}^{1}_{i,j}}_{a,j} = \operatorname{Min}_{k} \left\{ b + \mathcal{M}^{1}_{i,j-1}, c + \mathcal{M}'_{i,j} \right\} \end{cases}$$

Complexity :

For each min, $\mathcal{O}(n)$ potential contributors \Rightarrow Worst-case complexity in $\mathcal{O}(n^2)$ for naive backtrack. Keep best contributor for each Min \Rightarrow Backtracking in $\mathcal{O}(n)$

 \Rightarrow UnaFold [MZ08]/RNAFold [HFS⁺94] compute the MFE for the Turner model in **overall**¹ time/space complexities in $\mathcal{O}(n^3)/\mathcal{O}(n^2)$

^{1.} Using a trick/restriction for internal loops...

Definition (Ab initio folding)

Starting from sequence, find conformation that minimizes free-energy.

Advantages :

- Mechanical nature allows the (in)validation of models
- ► Reasonable complexity O(n³)/O(n²) time/space
- Exhaustive nature

Limitations :

- Hard to include PKs
- Highly dependent on energy model
- No cooperativity
- Limited performances

Definition (Comparative approach)

Starting from homologous sequences, postulate common structure and find best possible tradeoff between folding & alignment.

Avantages :

- Better performances
- (Limited) cooperativity
- Self-improving

Limitations

- Easily unreasonable complexity
- Non exhaustive search
- Captures transient structures

Definition (Ab initio folding)

Starting from sequence, find conformation that minimizes free-energy.

Advantages :

- Mechanical nature allows the (in)validation of models
- ► Reasonable complexity O(n³)/O(n²) time/space
- Exhaustive nature

Limitations :

- Hard to include PKs
- Highly dependent on energy model
- No cooperativity
- Limited performances

Definition (Comparative approach)

Starting from homologous sequences, postulate common structure and find best possible tradeoff between folding & alignment.

Avantages :

- Better performances
- (Limited) cooperativity
- Self-improving

Limitations

- Easily unreasonable complexity
- Non exhaustive search
- Captures transient structures







Outline

Focus on ab initio prediction

- Unambiguous Nussinov-like scheme
- Turner energy model
- Ab initio vs comparative

Boltzmann ensemble

- Nussinov : Minimisation ⇒ Comptage
- Calcul de la fonction de partition
- Échantillonnage statistique
- Inside/outside

Extensions

- Structures sous-optimales
- Pseudo-noeuds

L'ARN respire \Rightarrow II 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. ⇒ 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.

L'ARN respire \Rightarrow II 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. ⇒ 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.

Une distribution de Bolzmann pondère chaque structure *S* pour un ARN ω 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⁻¹)
- T est la température (K)
- ► *R* est la constante des gaz parfaits (1.986.10⁻³ kCal.K⁻¹.mol⁻¹)

Distribution renormalisée sur S_{ω} par la fonction de partition

$${\mathcal Z}_\omega = \sum_{{\mathcal S}\in {\mathcal S}_\omega} e^{rac{-{\mathcal E}_{{\mathcal S},\omega}}{RT}}.$$

où S_{ω} est l'ensemble des conformations compatibles avec ω .

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

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



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

$$N_{i,j} = \min \begin{cases} j & i \text{ unpaired} \\ \min_{k=i+\theta+1} \Delta G_{i,k} + N_{i+1,k-1} + N_{k+1,j} & i \text{ paired with } k \end{cases}$$

Ambiguity? Consider *i* : Either **unpaired**, or **paired** to *k*. Sets of structures generated in these two cases are clearly disjoint. (also holds for various values of k) \Rightarrow **Unambiguous** decomposition

Complete ? True, since scheme explores every possible outcome for *i*. + Induction on interval length \Rightarrow **Complete** decomposition



Récurrence sur l'énergie minimale d'un repliement :

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

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

Récurrence de comptage des structures compatibles :

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

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

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

Fonction de partition = Comptage pondéré des structures compatibles



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

$$\mathcal{Z}_{i,j} = \sum \begin{cases} \mathcal{Z}_{i+1,j} \\ \sum_{k=i+\theta+1}^{j} 1 \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \end{cases}$$


$$\begin{aligned} \mathcal{Z}_{i,t} &= 1, \quad \forall t \in [i, i + \theta] \\ \mathcal{Z}_{i,j} &= \sum \begin{cases} \mathcal{Z}_{i+1,j} \\ \sum_{k=i+\theta+1}^{j} e^{\frac{-\mathcal{E}_{bp}(i,k)}{AT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \end{aligned}$$



$$\begin{aligned} \mathcal{M}'_{i,j} &= \min \begin{cases} \frac{E_{\mathcal{H}}(i,j)}{E_{\mathcal{G}}(i,j) + \mathcal{M}'_{i+1,j-1}} \\ \min(E_{\mathcal{B}_{\ell}}(i,i',j',j) + \mathcal{M}'_{i',j'}) \\ a + c + \min(\mathcal{M}_{i+1,k-1} + \mathcal{M}^{1}_{k,j-1}) \end{cases} \\ \mathcal{M}_{i,j} &= \min \left\{ \min(\mathcal{M}_{i,k-1}, b(k-1)) + \mathcal{M}^{1}_{k,j} \right\} \\ \mathcal{M}^{1}_{i,j} &= \min \left\{ b + \mathcal{M}^{1}_{i,j-1}, c + \mathcal{M}'_{i,j} \right\} \end{aligned}$$



$$\mathcal{M}'_{i,j} = \operatorname{Min} \begin{cases} e^{\frac{-\mathcal{E}_{\mathcal{H}}^{(i,j)}}{\mathcal{H}^{\prime}}} + \mathcal{M}'_{i+1,j-1} \\ e^{\frac{-\mathcal{E}_{\mathcal{G}}^{(i,j)}}{\mathcal{H}^{\prime}}} + \mathcal{M}'_{i',j'} \\ e^{\frac{-(\theta+C)}{\mathcal{H}^{\prime}}} + \operatorname{Min}\left(\mathcal{M}_{i+1,k-1} + \mathcal{M}^{1}_{k,j-1}\right) \\ e^{\frac{-(\theta+C)}{\mathcal{H}^{\prime}}} + \operatorname{Min}\left(\mathcal{M}_{i,k-1}, e^{\frac{-\delta(k-1)}{\mathcal{H}^{\prime}}}\right) + \mathcal{M}^{1}_{k,j} \end{cases}$$
$$\mathcal{M}_{i,j} = \operatorname{Min} \left\{ e^{\frac{-\delta}{\mathcal{H}^{\prime}}} + \mathcal{M}^{1}_{i,j-1}, e^{\frac{-\delta(k-1)}{\mathcal{H}^{\prime}}} + \mathcal{M}'_{i,j} \right\}$$



$$\mathcal{M}'_{i,j} = \operatorname{Min} \begin{cases} e^{\frac{-\mathcal{E}_{H}(i,j)}{H^{\prime}}} \\ e^{\frac{-\mathcal{E}_{g}(i,j)}{H^{\prime}}} \mathcal{M}'_{i+1,j-1} \\ \operatorname{Min} \left(e^{\frac{-\mathcal{E}_{g}(i,i',j')}{H^{\prime}}} \mathcal{M}'_{i',j'} \right) \\ e^{\frac{-(a+c)}{H^{\prime}}} \operatorname{Min} \left(\mathcal{M}_{i+1,k-1} \mathcal{M}^{1}_{k,j-1} \right) \\ \mathcal{M}_{i,j} = \operatorname{Min} \left\{ \operatorname{Min} \left(\mathcal{M}_{i,k-1}, e^{\frac{-\partial(k-1)}{H^{\prime}}} \right) \mathcal{M}^{1}_{k,j} \right\} \\ \mathcal{M}^{1}_{i,j} = \operatorname{Min} \left\{ e^{\frac{-b}{H^{\prime}}} \mathcal{M}^{1}_{i,j-1}, e^{\frac{-c}{H^{\prime}}} \mathcal{M}'_{i,j} \right\} \end{cases}$$



$$\begin{aligned} \mathcal{Z}'(i,j) &= \sum \begin{cases} e^{\frac{-E_{H}(i,j)}{H^{1}}} \\ e^{\frac{-E_{H}(i,j)}{H^{1}}} \mathcal{Z}'(i+1,j-1) \\ + \sum \left(e^{\frac{-E_{H}(i,j',j',j)}{H^{1}}} \mathcal{Z}'(i',j') \right) \\ + e^{\frac{-(b+c)}{H^{1}}} \sum \left(\mathcal{Z}(i+1,k-1)\mathcal{Z}^{1}(k,j-1) \right) \end{aligned}$$
$$\begin{aligned} \mathcal{Z}(i,j) &= \sum \left(\mathcal{Z}(i,k-1) + e^{\frac{-b(k-1)}{H^{1}}} \right) \mathcal{Z}^{1}(k,j) \\ \mathcal{Z}^{1}(i,j) &= e^{\frac{-b}{H}} \mathcal{Z}^{1}(i,j-1) + e^{\frac{-c}{H}} \mathcal{Z}'(i,j) \end{aligned}$$

Fonction de partition = Comptage pondéré des structures compatibles

$$\begin{aligned} \mathcal{Z}_{i,t} &= 1, \quad \forall t \in [i, i + \theta] \\ \mathcal{Z}_{i,j} &= \sum \begin{cases} \mathcal{Z}_{i+1,j} \\ \sum_{k=i+\theta+1}^{j} e^{\frac{-E_{\text{bp}}(i,k)}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \end{aligned}$$

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

$$\left(e^{-a/RT} \cdot Z^1 \cdot Z' = e^{-a/RT} \cdot \sum_x e^{-E_x/RT} \cdot \sum_y e^{-E_y/RT}\right)$$

 $= \sum_{x,y} e^{-a/RT} \cdot e^{-E_x/RT} \cdot e^{-E_y/RT} = \sum_{x,y} e^{-(a+E_x+E_y)/RT}$

Fonction de partition = Comptage pondéré des structures compatibles

$$\begin{aligned} \mathcal{Z}_{i,t} &= 1, \quad \forall t \in [i, i+\theta] \\ \mathcal{Z}_{i,j} &= \sum \begin{cases} \mathcal{Z}_{i+1,j} \\ \sum_{k=i+\theta+1}^{j} e^{-\frac{E_{bp}(i,k)}{HT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \end{aligned}$$

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 $\left(e^{-a/RT} \cdot Z^1 \cdot Z' = e^{-a/RT} \cdot \sum_x e^{-E_x/RT} \cdot \sum_y e^{-E_y/RT}\right)$ $= \sum_{x,y} e^{-a/RT} \cdot e^{-E_x/RT} \cdot e^{-E_y/RT} = \sum_{x,y} e^{-(a+E_x+E_y)/RT}$



Fonction de partition = Comptage pondéré des structures compatibles

$$\begin{aligned} \mathcal{Z}_{i,t} &= 1, \quad \forall t \in [i, i+\theta] \\ \mathcal{Z}_{i,j} &= \sum \begin{cases} \mathcal{Z}_{i+1,j} \\ \sum_{k=i+\theta+1}^{j} e^{-\frac{E_{bp}(i,k)}{HT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \end{aligned}$$

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 $\left(e^{-a/RT} \cdot \mathbf{Z}^{1} \cdot \mathbf{Z}' = e^{-a/RT} \cdot \sum_{x} e^{-E_{x}/RT} \cdot \sum_{y} e^{-E_{y}/RT}$ $= \sum_{x,y} e^{-a/RT} \cdot e^{-E_{x}/RT} \cdot e^{-E_{y}/RT} = \sum_{x,y} e^{-(a+E_{x}+E_{y})/RT}\right)$



Fonction de partition = Comptage pondéré des structures compatibles

$$\begin{aligned} \mathcal{Z}_{i,t} &= 1, \quad \forall t \in [i, i+\theta] \\ \mathcal{Z}_{i,j} &= \sum \begin{cases} \mathcal{Z}_{i+1,j} \\ \sum_{k=i+\theta+1}^{j} e^{\frac{-E_{bp}(i,k)}{HT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \end{aligned}$$

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 $\left(e^{-a/RT} \cdot \mathbf{Z}^{1} \cdot \mathbf{Z}' = e^{-a/RT} \cdot \sum_{x} e^{-E_{x}/RT} \cdot \sum_{y} e^{-E_{y}/RT}$ $= \sum_{x,y} e^{-a/RT} \cdot e^{-E_{x}/RT} \cdot e^{-E_{y}/RT} = \sum_{x,y} e^{-(a+E_{x}+E_{y})/RT}\right)$



Fonction de partition = Comptage pondéré des structures compatibles

$$\begin{aligned} \mathcal{Z}_{i,t} &= 1, \quad \forall t \in [i, i+\theta] \\ \mathcal{Z}_{i,j} &= \sum \begin{cases} \mathcal{Z}_{i+1,j} \\ \sum_{k=i+\theta+1}^{j} e^{\frac{-E_{bp}(i,k)}{HT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \end{aligned}$$

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 $\left(e^{-a/RT} \cdot \mathbf{Z}^{1} \cdot \mathbf{Z}' = e^{-a/RT} \cdot \sum_{x} e^{-E_{x}/RT} \cdot \sum_{y} e^{-E_{y}/RT}$ $= \sum_{x,y} e^{-a/RT} \cdot e^{-E_{x}/RT} \cdot e^{-E_{y}/RT} = \sum_{x,y} e^{-(a+E_{x}+E_{y})/RT}\right)$



Fonction de partition = Comptage pondéré des structures compatibles

$$\begin{aligned} \mathcal{Z}_{i,t} &= 1, \quad \forall t \in [i, i+\theta] \\ \mathcal{Z}_{i,j} &= \sum \begin{cases} \mathcal{Z}_{i+1,j} \\ \sum_{k=i+\theta+1}^{j} e^{\frac{-E_{bp}(i,k)}{HT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \end{aligned}$$

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 $\left(e^{-a/RT} \cdot Z^1 \cdot Z' = e^{-a/RT} \cdot \sum_x e^{-E_x/RT} \cdot \sum_y e^{-E_y/RT}\right)$ $= \sum_{x,y} e^{-a/RT} \cdot e^{-E_x/RT} \cdot e^{-E_y/RT} = \sum_{x,y} e^{-(a+E_x+E_y)/RT}$



Fonction de partition = Comptage pondéré des structures compatibles

$$\begin{aligned} \mathcal{Z}_{i,t} &= 1, \quad \forall t \in [i, i+\theta] \\ \mathcal{Z}_{i,j} &= \sum \begin{cases} \mathcal{Z}_{i+1,j} \\ \sum_{k=i+\theta+1}^{j} e^{-\frac{E_{bp}(i,k)}{HT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \end{aligned}$$

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 $\left(e^{-a/RT} \cdot Z^1 \cdot Z' = e^{-a/RT} \cdot \sum_x e^{-E_x/RT} \cdot \sum_y e^{-E_y/RT}\right)$ $= \sum_{x,y} e^{-a/RT} \cdot e^{-E_x/RT} \cdot e^{-E_y/RT} = \sum_{x,y} e^{-(a+E_x+E_y)/RT}$



Fonction de partition = Comptage pondéré des structures compatibles

$$\begin{aligned} \mathcal{Z}_{i,t} &= 1, \quad \forall t \in [i, i+\theta] \\ \mathcal{Z}_{i,j} &= \sum \begin{cases} \mathcal{Z}_{i+1,j} \\ \sum_{k=i+\theta+1}^{j} e^{-\frac{E_{bp}(i,k)}{HT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \end{aligned}$$

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 $\left(e^{-a/RT} \cdot Z^1 \cdot Z' = e^{-a/RT} \cdot \sum_x e^{-E_x/RT} \cdot \sum_y e^{-E_y/RT}\right)$ $= \sum_{x,y} e^{-a/RT} \cdot e^{-E_x/RT} \cdot e^{-E_y/RT} = \sum_{x,y} e^{-(a+E_x+E_y)/RT}$

Exemple :



Yann Ponty (CNRS, Ecole Polytechnique) AMI2B Master - ABA Lecture - Boltzmann ensemble

Fonction de partition = Comptage pondéré des structures compatibles

$$\begin{aligned} \mathcal{Z}_{i,t} &= 1, \quad \forall t \in [i, i+\theta] \\ \mathcal{Z}_{i,j} &= \sum \begin{cases} \mathcal{Z}_{i+1,j} \\ \sum_{k=i+\theta+1}^{j} e^{-\frac{E_{bp}(i,k)}{HT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \end{aligned}$$

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 $\left(e^{-a/RT} \cdot Z^1 \cdot Z' = e^{-a/RT} \cdot \sum_x e^{-E_x/RT} \cdot \sum_y e^{-E_y/RT}\right)$ $= \sum_{x,y} e^{-a/RT} \cdot e^{-E_x/RT} \cdot e^{-E_y/RT} = \sum_{x,y} e^{-(a+E_x+E_y)/RT}$



Fonction de partition = Comptage pondéré des structures compatibles

$$\begin{aligned} \mathcal{Z}_{i,t} &= 1, \quad \forall t \in [i, i+\theta] \\ \mathcal{Z}_{i,j} &= \sum \begin{cases} \mathcal{Z}_{i+1,j} \\ \sum_{k=i+\theta+1}^{j} e^{-\frac{E_{bp}(i,k)}{HT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \end{aligned}$$

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 $\left(e^{-a/RT} \cdot Z^1 \cdot Z' = e^{-a/RT} \cdot \sum_x e^{-E_x/RT} \cdot \sum_y e^{-E_y/RT}\right)$ $= \sum_{x,y} e^{-a/RT} \cdot e^{-E_x/RT} \cdot e^{-E_y/RT} = \sum_{x,y} e^{-(a+E_x+E_y)/RT}$



Fonction de partition = Comptage pondéré des structures compatibles

$$\begin{aligned} \mathcal{Z}_{i,t} &= 1, \quad \forall t \in [i, i+\theta] \\ \mathcal{Z}_{i,j} &= \sum \begin{cases} \mathcal{Z}_{i+1,j} \\ \sum_{k=i+\theta+1}^{j} e^{-\frac{E_{bp}(i,k)}{HT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \end{aligned}$$

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 $\left(e^{-a/RT} \cdot Z^1 \cdot Z' = e^{-a/RT} \cdot \sum_x e^{-E_x/RT} \cdot \sum_y e^{-E_y/RT}\right)$ $= \sum_{x,y} e^{-a/RT} \cdot e^{-E_x/RT} \cdot e^{-E_y/RT} = \sum_{x,y} e^{-(a+E_x+E_y)/RT}$



Fonction de partition = Comptage pondéré des structures compatibles

$$\begin{aligned} \mathcal{Z}_{i,t} &= 1, \quad \forall t \in [i, i+\theta] \\ \mathcal{Z}_{i,j} &= \sum \begin{cases} \mathcal{Z}_{i+1,j} \\ \sum_{k=i+\theta+1}^{j} e^{-\frac{E_{bp}(i,k)}{HT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \end{aligned}$$

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 $\left(e^{-a/RT} \cdot Z^1 \cdot Z' = e^{-a/RT} \cdot \sum_x e^{-E_x/RT} \cdot \sum_y e^{-E_y/RT}\right)$ $= \sum_{x,y} e^{-a/RT} \cdot e^{-E_x/RT} \cdot e^{-E_y/RT} = \sum_{x,y} e^{-(a+E_x+E_y)/RT}$



Echantillonnage statistique de structures d'ARN

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

 \Rightarrow 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
- \Rightarrow 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?

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

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

8 Réitérer sur les sous-structures

$$\mathcal{Z}'(i,j) \in \left\{ \begin{array}{c} - - \partial e^{\frac{-E_{H}(i,j)}{RT}} + e^{\frac{-E_{S}(i,j)}{RT}} \mathcal{Z}'(i+1,j-1) \\ \oplus e^{\frac{-E_{B}(i,i',j',j)}{RT}} \mathcal{Z}'(i',j') \end{array} \right\}$$

$$B$$

$$B$$

$$B$$

$$B$$

$$B$$

$$C$$

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

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

2) Retirer à *r* les contributions à $\mathcal{Z}'(i, j)$, jusqu'à ce que r < 0

8 Réitérer sur les sous-structures

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

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

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

2) Retirer à *r* les contributions à $\mathcal{Z}'(i, j)$, jusqu'à ce que r < 0

8 Réitérer sur les sous-structures

$$\mathcal{Z}'(i,j) = \sum \begin{cases} e^{\frac{-E_{H}(i,j)}{RT}} + e^{\frac{-E_{S}(i,j)}{RT}} \mathcal{Z}'(i+1,j-1) & \mathbb{A} \\ \sum \left(e^{\frac{-E_{BI}(i,i',j',j)}{RT}} \mathcal{Z}'(i',j') \right) & \mathbb{B} \\ e^{\frac{-(a+c)}{RT}} \sum \left(\mathcal{Z}(i+1,k-1) \mathcal{Z}^{1}(k,j-1) \right) & \mathbb{C} \\ & \downarrow \\ \mathbf{A}_{1} \| \mathbf{A}_{2} \| \mathbf{B}_{j} \| \mathbf{B}_{i+1} \| \dots \| \mathbf{B}_{j-1} \| \mathbf{B}_{j} \| \mathbf{C}_{i} \| \mathbf{C}_{i+1} \| \dots \| \mathbf{C}_{j-1} \| \mathbf{C}_{j} \end{cases}$$

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

- 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

8 Réitérer sur les sous-structures

$$\mathcal{Z}'(i,j) = \sum \begin{cases} e^{\frac{-E_{H}(i,j)}{RT}} + e^{\frac{-E_{S}(i,j)}{RT}} \mathcal{Z}'(i+1,j-1) & \mathbb{A} \\ \sum \left(e^{\frac{-E_{BI}(i,i',j',j)}{RT}} \mathcal{Z}'(i',j') \right) & \mathbb{B} \\ e^{\frac{-(a+c)}{RT}} \sum \left(\mathcal{Z}(i+1,k-1)\mathcal{Z}^{1}(k,j-1) \right) & \mathbb{C} \\ & \downarrow \\ & \downarrow \\ \mathbf{A}_{1} | \mathbf{A}_{2} | \mathbf{B}_{j} | \mathbf{B}_{i+1} | \dots | \mathbf{B}_{j-1} | \mathbf{B}_{j} | \mathbf{C}_{i} | \mathbf{C}_{i+1} | \dots | \mathbf{C}_{j-1} | \mathbf{C}_{j} \end{cases}$$

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

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

8 Réitérer sur les sous-structures

$$\mathcal{Z}'(i,j) = \sum \begin{cases} e^{\frac{-E_{H}(i,j)}{RT}} + e^{\frac{-E_{S}(i,j)}{RT}} \mathcal{Z}'(i+1,j-1) & \mathbb{A} \\ \sum \left(e^{\frac{-E_{BI}(i,i',j',j)}{RT}} \mathcal{Z}'(i',j') \right) & \mathbb{B} \\ e^{\frac{-(a+c)}{RT}} \sum \left(\mathcal{Z}(i+1,k-1)\mathcal{Z}^{1}(k,j-1) \right) & \mathbb{C} \\ & \downarrow \\ \mathbf{A}_{1} | \mathbf{A}_{2} | \mathbf{B}_{i} | \mathbf{B}_{i+1} | \dots | \mathbf{B}_{j-1} | \mathbf{B}_{j} | \mathbf{C}_{i} | \mathbf{C}_{i+1} | \dots | \mathbf{C}_{j-1} | \mathbf{C}_{j} \end{cases}$$

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

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

8 Réitérer sur les sous-structures

$$\mathcal{Z}'(i,j) = \sum \begin{cases} e^{\frac{-E_{H}(i,j)}{RT}} + e^{\frac{-E_{S}(i,j)}{RT}} \mathcal{Z}'(i+1,j-1) & \mathbb{A} \\ \sum \left(e^{\frac{-E_{BI}(i,i',j',j)}{RT}} \mathcal{Z}'(i',j') \right) & \mathbb{B} \\ e^{\frac{-(a+c)}{RT}} \sum \left(\mathcal{Z}(i+1,k-1)\mathcal{Z}^{1}(k,j-1) \right) & \mathbb{C} \\ & \downarrow \\ &$$

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

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

8 Réitérer sur les sous-structures

$$\mathcal{Z}'(i,j) = \sum \begin{cases} e^{\frac{-E_{H}(i,j)}{RT}} + e^{\frac{-E_{S}(i,j)}{RT}} \mathcal{Z}'(i+1,j-1) & \mathbb{A} \\ \sum \left(e^{\frac{-E_{BI}(i,i',j',j)}{RT}} \mathcal{Z}'(i',j') \right) & \mathbb{B} \\ e^{\frac{-(a+c)}{RT}} \sum \left(\mathcal{Z}(i+1,k-1)\mathcal{Z}^{1}(k,j-1) \right) & \mathbb{C} \\ & \downarrow \\ \mathbf{A}_{1} | \mathbf{A}_{2} | \mathbf{B}_{j} | \mathbf{B}_{i+1} | \dots | \mathbf{B}_{j-1} | \mathbf{B}_{j} | \mathbf{C}_{i} | \mathbf{C}_{i+1} | \dots | \mathbf{C}_{j-1} | \mathbf{C}_{j} \\ & \downarrow \\ \mathbf{A}_{1} | \mathbf{A}_{2} | \mathbf{B}_{j} | \mathbf{B}_{i+1} | \dots | \mathbf{B}_{j-1} | \mathbf{B}_{j} | \mathbf{C}_{i} | \mathbf{C}_{i+1} | \dots | \mathbf{C}_{j-1} | \mathbf{C}_{j} \end{cases}$$

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

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

$$\mathcal{Z}'(i,j) = \sum \begin{cases} e^{-\frac{E_{H}(i,j)}{RT}} + e^{-\frac{E_{S}(i,j)}{RT}} \mathcal{Z}'(i+1,j-1) & \mathbb{A} \\ \sum \left(e^{-\frac{E_{B}(i,i',j',j)}{RT}} \mathcal{Z}'(i',j') \right) & \mathbb{B} \end{cases}$$

Correction : Tout $S \in S_{\omega}$ engendrée de façon unique (Unambiguité de Turner)

La probabilité d'engendrer S est donc

$$\boldsymbol{\rho}_{\mathcal{S}} = \frac{\mathcal{B}(E_1)}{\mathcal{B}(\mathcal{S}_{\mathbf{w}})} \cdot \frac{\mathcal{B}(E_2)}{\mathcal{B}(E_1)} \cdot \frac{\mathcal{B}(E_3)}{\mathcal{B}(E_2)} \cdots \frac{\mathcal{B}(\{\boldsymbol{S}\})}{\mathcal{B}(E_m)}$$

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

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

$$\mathcal{Z}'(i,j) = \sum \begin{cases} e^{\frac{-\mathcal{E}_{\mathcal{H}}(i,j)}{RT}} + e^{\frac{-\mathcal{E}_{\mathcal{S}}(i,j)}{RT}} \mathcal{Z}'(i+1,j-1) & \\ \sum \left(e^{\frac{-\mathcal{E}_{\mathcal{B}}(i,i',j',j)}{RT}} \mathcal{Z}'(i',j') \right) & \\ \end{cases}$$

$$e^{\frac{-(a+c)}{RT}}\sum \left(\mathcal{Z}(i+1,k-1)\mathcal{Z}^{1}(k,j-1)\right) \quad \mathbb{C}$$

Correction : Tout $S \in S_{\omega}$ engendrée de façon unique (Unambiguité de Turner)

La probabilité d'engendrer S est donc

$$p_{S} = \frac{1}{\mathcal{B}(\mathcal{S}_{W})} \cdot \frac{1}{1} \cdot \frac{1}{1} \dots \frac{\mathcal{B}(\{S\})}{1}$$

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

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

$$e^{\frac{-E_{H}(i,j)}{RT}} + e^{\frac{-E_{S}(i,j)}{RT}} \mathcal{Z}'(i+1,j-1) \qquad (A)$$

$$Z'(i,j) = \sum \left\{ \sum \left(e^{-Rt} Z'(r,j') \right) e^{\frac{-(a+c)}{RT}} \sum \left(Z(i+1,k-1)Z^{1}(k,j-1) \right) \right\}$$

Correction : Tout $S \in S_{\omega}$ engendrée de façon unique (Unambiguité de Turner)

La probabilité d'engendrer S est donc

$$p_{S} = \frac{\mathcal{B}(\{S\})}{\mathcal{B}(\mathcal{S}_{W})} = \frac{e^{-E_{S}/RT}}{\mathcal{Z}} = P_{S,\omega}$$
Complexité ? ? ?

Complexité

Algorithme (Reformulation SFold [DL03a])

- 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
- 8 Réitérer sur les sous-structures

$$Z'(i,j) \in \underbrace{\left\{\begin{array}{c} -- \Rightarrow e^{\frac{-E_{H}(i,j)}{RT}} + e^{\frac{-E_{S}(i,j)}{RT}} Z'(i+1,j-1) \\ \Rightarrow e^{\frac{-(i+c)}{RT}} \\ \Rightarrow e^{\frac{-(a+c)}{RT}} \sum \left(Z(i+1,k-1) Z^{1}(k,j-1) \right) \\ \end{array} \right\}} \begin{bmatrix} e^{-\frac{(a+c)}{RT}} \sum \left(Z(i+1,k-1) Z^{1}(k,j-1) \right) \\ \Rightarrow e^{\frac{-(a+c)}{RT}} \sum \left(Z(i+1,k-1) Z^{1}(k,j-1) \right) \\ \Rightarrow e^{\frac{-(a+c)}{RT}} \\ \Rightarrow e^{\frac{-(a+c)}{RT}} \sum \left(Z(i+1,k-1) Z^{1}(k,j-1) \right) \\ \Rightarrow e^{\frac{-(a+c)}{RT}} \\ \Rightarrow e^{\frac{-(a+c)}{RT}} \sum \left(Z(i+1,k-1) Z^{1}(k,j-1) \right) \\ \Rightarrow e^{\frac{-(a+c)}{RT}} \\ \Rightarrow e^{\frac{-(a$$

Complexité en moyenne en $\Theta(n\sqrt{n})$ dans l'hypothèse tout appariement. Adaptation d'un parcours Boustrophedon $\Rightarrow O(n \log nk)$ au pire.

Complexité

Algorithme (Reformulation SFold [DL03a])

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

$$Z'(i,j) = \sum \begin{cases} e^{-\frac{E_{H}(i,j)}{RT}} + e^{-\frac{E_{S}(i,j)}{RT}} Z'(i+1,j-1) & A \\ \sum \left(e^{-\frac{E_{B}(i,i',j',j)}{RT}} Z'(i',j') \right) & B \\ e^{-\frac{(a+c)}{RT}} \sum \left(Z(i+1,k-1) Z^{1}(k,j-1) \right) & C \\ \downarrow & \downarrow \\ A_{1} | A_{2} | B_{i} | B_{i+1} | \dots | B_{j-1} | B_{j} | C_{j} | C_{i+1} | \dots | C_{j-1} | C_{j} \\ \downarrow & \downarrow \\ C \\ A_{1} | A_{2} | B_{i} | B_{i+1} | \dots | B_{j-1} | B_{j} | C_{j} | C_{i+1} | \dots | C_{j-1} | C_{j} \end{cases}$$

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

Complexité en moyenne en $\Theta(n\sqrt{n})$ dans l'hypothèse tout appariement. Adaptation d'un parcours Boustrophedon $\Rightarrow O(n \log nk)$ au pire.
















Si le schéma de prog. dyn. est **acyclique** et **indépendant**, cette décomposition est **complète** et **non-ambiguë**, et implique une récurrence *simple* pour le probabilités de paires de bases ...

Outline

Focus on ab initio prediction

- Unambiguous Nussinov-like scheme
- Turner energy model
- Ab initio vs comparative

2 Boltzmann ensemble

- Nussinov : Minimisation ⇒ Comptage
- Calcul de la fonction de partition
- Échantillonnage statistique
- Inside/outside

Extensions

- Structures sous-optimales
- Pseudo-noeuds

- Prob. : Simplifications de l'énergie (Pseudo-noeuds, non-can.)
- \Rightarrow La structure **native** (fonctionnelle) pourrait être **ignorée**.

- 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 \geq 1 struct.
- Engendrer (Rec.) les sous-ensembles et combiner (brutal ou Tri)

$$\mathcal{M}'_{1,n,\Delta} = \varepsilon = (-----)^{-1} (a+c+\operatorname{Min}\left(\mathcal{M}_{i+1,k_0-1}+\mathcal{M}^{1}_{k_0,j-1}\right)) \qquad E_0 - \mathcal{M}'_{1,n} = \varepsilon_0 \leq \Delta$$

$$\mathcal{M}'_{1,n,\Delta} = \varepsilon = (------)^{-1} (a+c+\operatorname{Min}\left(\mathcal{M}_{i+1,k_1-1}+\mathcal{M}^{1}_{k_1,j-1}\right)) \qquad E_1 - \mathcal{M}'_{1,n} = \varepsilon_1 > \Delta$$

$$\mathcal{M}'_{1,n,\Delta} = \varepsilon = (-------)^{-1} (a+c+\operatorname{Min}\left(\mathcal{M}_{i+1,k_2-1}+\mathcal{M}^{1}_{k_2,j-1}\right)) \qquad E_2 - \mathcal{M}'_{1,n} = \varepsilon_2 \leq \Delta$$

- Prob. : Simplifications de l'énergie (Pseudo-noeuds, non-can.)
- \Rightarrow La structure **native** (fonctionnelle) pourrait être **ignorée**.

- 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 \geq 1 struct.
- Engendrer (Rec.) les sous-ensembles et combiner (brutal ou Tri)

$$\begin{array}{c}
\mathcal{M}'_{1,n,\Delta} \\
\mathcal{M}'_{1,n,\Delta} \\
\mathcal{M}'_{k_0,j-1} \\
\Delta' = \Delta - \varepsilon_0
\end{array}$$

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

- 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 \geq 1 struct.
- ► Engendrer (Rec.) les sous-ensembles et combiner (brutal ou Tri)



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

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



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

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



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

 \Rightarrow 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 \geq 1 struct.
- Engendrer (Rec.) les sous-ensembles et combiner (brutal ou Tri)

 $\Rightarrow \text{Complexité en temps (Tri)} : \mathcal{O}(n^3 + nk\log(k))$ (k croît exponentiellement sur Δ !)

Repliement avec pseudo-noeuds

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



Leur absence historique au sein algorithmes de repliement est liée à la difficulté algorithmique des problèmes associés.

(Présence de croisements interdit une hypothèse d'indépendance des repliements).

Туре	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]

But : Capturer des catégorie de pseudo-noeuds *simples*, mais très représentés.



ldée : Quand on *retourne* ce type de pseudonoeuds, il suffit de précalculer les meilleures configurations *en dessous* d'un triplet (i, j, k) pour obtenir son énergie minimale.

But: Capturer des catégorie de pseudo-noeuds *simples*, mais très représentés.



Idée : Quand on *retourne* ce type de pseudonoeuds, il suffit de précalculer les meilleures configurations *en dessous* d'un triplet (i, j, k) pour obtenir son énergie minimale.

But: Capturer des catégorie de pseudo-noeuds *simples*, mais très représentés.



Idée : Quand on *retourne* ce type de pseudonoeuds, il suffit de précalculer les meilleures configurations *en dessous* d'un triplet (i, j, k) pour obtenir son énergie minimale.

Akutsu/Uemura : Programmation dynamique



Application/Problème	Weight fun.	Time/Space	Ref.
Minimisation d'énergie	π_{bp}	$\mathcal{O}(n^4)/\mathcal{O}(n^4)$	[Aku00]
Fonction de partition	$e^{\frac{-\pi_{bp}}{RT}}$	$\mathcal{O}(n^4)/\mathcal{O}(n^4)$	$\Theta(n^6)$ [CC09]
Probabilité de paires de bases	e ^{-m} _{bp}	$\mathcal{O}(n^4)/\mathcal{O}(n^4)$	-
Échantillonnage (k-struct.)	e ^{-m} _{bp}	$\mathcal{O}(n^4 + kn \log n) / \mathcal{O}(n^4)$	-

Exercice : Ecrire l'équation de programmation dynamique associée pour le repliement, le comptage et la fonction de partition.

References I



Tatsuya Akutsu.

Dynamic programming algorithms for rna secondary structure prediction with pseudoknots. Discrete Appl. Math., 104(1-3):45–62, 2000.



S. Cao and S-J Chen.

Predicting structured and stabilities for h-type pseudoknots with interhelix loop. RNA, 15 :696–706, 2009.



K. Doshi, J. J. Cannone, C. Cobaugh, and R. R. Gutell.

Evaluation of the suitability of free-energy minimization using nearest-neighbor energy parameters for rna secondary structure prediction.

BMC Bioinformatics, 5(1) :105, 2004.



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.



Y. Ding and E. Lawrence.

A statistical sampling algorithm for RNA secondary structure prediction. *Nucleic Acids Res*, 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.



P. Gardner and R. Giegerich.

A comprehensive comparison of comparative rna structure prediction approaches. BMC Bioinformatics, 5(1) :140, 2004.

References II



I. L. Hofacker, W. Fontana, P. F. Stadler, L. S. Bonhoeffer, M. Tacker, and P. Schuster.

Fast folding and comparison of RNA secondary structures. Monatshefte für Chemie / Chemical Monthly, 125(2) :167–188, 1994.



R. B. Lyngsøand C. N. S. Pedersen.

RNA pseudoknot prediction in energy-based models. Journal of Computational Biology, 7(3-4) :409–427, 2000.



J.S. McCaskill.

The equilibrium partition function and base pair binding probabilities for RNA secondary structure. *Biopolymers*, 29:1105–1119, 1990.

István Miklós, Irmtraud M Meyer, and Borbála Nagy.

Moments of the boltzmann distribution for rna secondary structures. *Bull Math Biol*, 67(5) :1031–1047, Sep 2005.



Expanded sequence dependence of thermodynamic parameters improves prediction of RNA secondary structure. J Mol Biol, 288 :911–940, 1999.



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. J Math Biol, 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 :2053–2068, 1999.



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.



M. Zuker and P. Stiegler.

Optimal computer folding of large RNA sequences using thermodynamics and auxiliary information. *Nucleic Acids Res*, 9:133–148, 1981.