

Algorithmics and combinatorics of RNA sampling

Danièle Gardi[‡] Andy Lorenz[†] Yann Ponty^{*}

[‡] Université Versailles St Quentin – France

[†] Boston College – Boston – USA

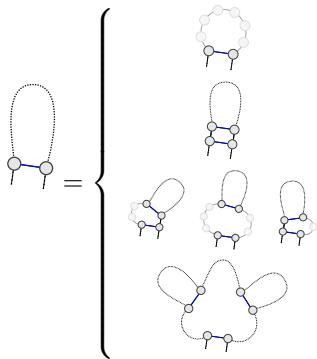
^{*} Polytechnique/CNRS/INRIA AMIB – France

July 29, 2009

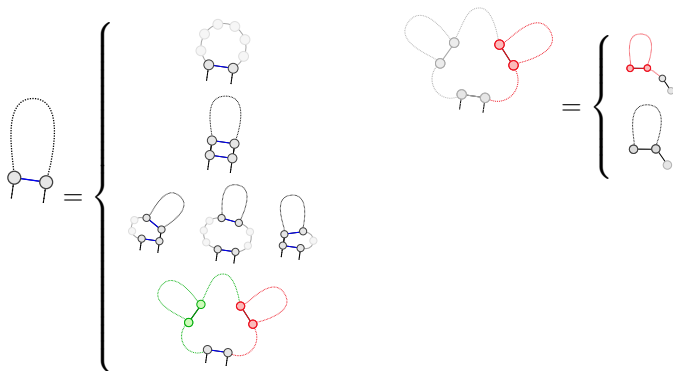
(Enumerative) Combinatorics helps:

- Counting conformations
- Analyzing features of null models
- Performing analysis of algorithms. . .
- . . . improving them?

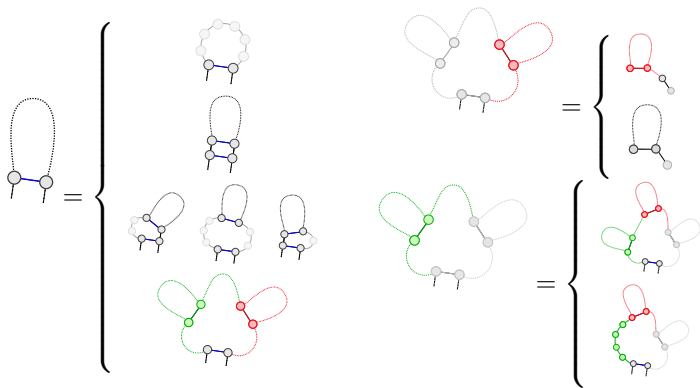
MFE DP equations



MFE DP equations



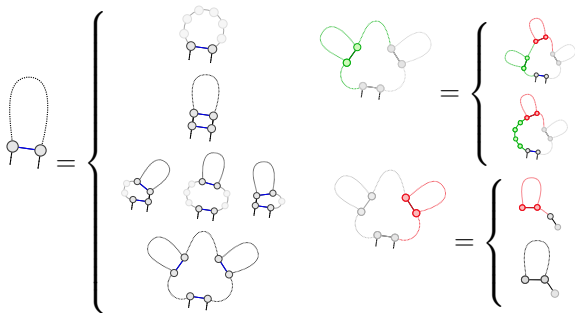
MFE DP equations



Completeness: Not too hard to check!

Unambiguity \Rightarrow Enumerative combinatorics.

Validation



Generating function $T(z) = \sum_{n \geq 0} t_n z^n$

With $t_n = \#$ Secondary structures of size n

$$\mathcal{A}(z) = \begin{cases} S(z) \\ z^2 \mathcal{A}(z) \\ zS(z)z^2 \mathcal{A}(z) + z^2 \mathcal{A}(z)S(z)z \\ + zS(z)z^2 \mathcal{A}(z)S(z)z \\ B(z)\mathcal{C}(z) \end{cases} \quad \begin{cases} \mathcal{B}(z) = \begin{cases} B(z)\mathcal{C}(z) \\ S(z)\mathcal{B}(z) \end{cases} \\ \mathcal{C}(z) = \begin{cases} \mathcal{C}(z)z \\ z^2 \mathcal{A}(z) \end{cases} \end{cases}$$

$$S(z) = 1 + zS(z)$$

$$A(z) = \begin{cases} S(z) \\ z^2 A(z) \\ zS(z)z^2 A(z) + z^2 A(z)S(z)z \\ + zS(z)z^2 A(z)S(z)z \\ B(z)C(z) \end{cases} \quad \begin{cases} B(z) = \begin{cases} B(z)C(z) \\ S(z)B(z) \end{cases} \\ C(z) = \begin{cases} C(z)z \\ z^2 A(z) \end{cases} \end{cases}$$

$$S(z) = 1 + zS(z)$$

Reminder: Waterman counted Sec. Str. [Wat78] and found the gen. fun.

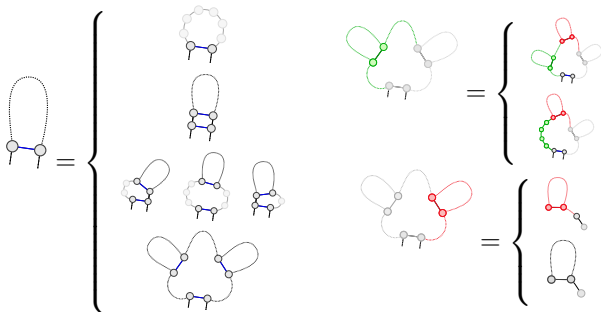
$$\mathcal{W}(z) = \frac{1 - z + z^2 - \sqrt{1 - 2z - z^2 - 2z^3 + z^4}}{2z^2}$$

Here we have

$$\begin{aligned} \Rightarrow A(z) &= \frac{1 - z - z^2 - \sqrt{1 - 2z - z^2 - 2z^3 + z^4}}{2z^2} \\ &= \mathcal{W}(z) - 1 \quad (\text{Woops, we forgot the empty RNA}) \end{aligned}$$

MFE folding

- $E_H(i, j)$: Energy of hairpin loop with closing pair (i, j)
- $E_{BI}(i, j)$: Energy of bulge or internal loop with closing pair (i, j)
- $E_S(i, j)$: Energy of stacking pairs $(i, j)/(i + 1, j - 1)$
- a, c, b : Penalties for multiloop, hairpins and unpaired bases in multiloop.



Message #1

Treating search space as a combinatorial object saves time and trouble!

MFE folding

- $E_H(i,j)$: Energy of hairpin loop with closing pair (i,j)
- $E_{BI}(i,j)$: Energy of bulge or internal loop with closing pair (i,j)
- $E_S(i,j)$: Energy of stacking pairs $(i,j)/(i+1,j-1)$
- a,c,b : Penalties for multiloop, hairpins and unpaired bases in multiloop.

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

Message #1

Treating search space as a combinatorial object saves time and trouble!

Partition function/Boltzmann probability

- Let ω be an RNA sequence
- \mathcal{S}_ω be the set of sequences compatible with ω ,

$$\text{Partition function } Z_\omega = \sum_{S \in \mathcal{S}_\omega} e^{\frac{-E_{S,\omega}}{RT}}$$

where T is temperature in Kelvin and R is the universal gas constant.

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

$P_{S,\omega}$ is the probability of observing ω in conformation S .

- ⇒ Offers a more dynamic view of the folding process
- ⇒ Gives a model for computing various probabilities (BP, Motifs ...)
- ⇒ Unified algorithmic framework for subopts and mfe ($RT \rightarrow \infty$)
- ⇒ Very easy to embed into any existing DP equations

Partition function: Beyond m.f.e. hypothesis

Partition function/Boltzmann probability

- Let ω be an RNA sequence
- \mathcal{S}_ω be the set of sequences compatible with ω ,

$$\text{Partition function } Z_\omega = \sum_{S \in \mathcal{S}_\omega} e^{\frac{-E_{S,\omega}}{RT}}$$

where T is temperature in Kelvin and R is the universal gas constant.

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

$P_{S,\omega}$ is the probability of observing ω in conformation S .

- ⇒ Offers a more dynamic view of the folding process
- ⇒ Gives a model for computing various probabilities (BP, Motifs ...)
- ⇒ Unified algorithmic framework for subopts and mfe ($RT \rightarrow \infty$)
- ⇒ Very easy to embed into any existing DP equations

Partition function: Beyond m.f.e. hypothesis

Partition function/Boltzmann probability

- Let ω be an RNA sequence
- \mathcal{S}_ω be the set of sequences compatible with ω ,

$$\text{Partition function } Z_\omega = \sum_{S \in \mathcal{S}_\omega} e^{\frac{-E_{S,\omega}}{RT}}$$

where T is temperature in Kelvin and R is the universal gas constant.

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

$P_{S,\omega}$ is the probability of observing ω in conformation S .

- ⇒ Offers a more dynamic view of the folding process
- ⇒ Gives a model for computing various probabilities (BP, Motifs ...)
- ⇒ Unified algorithmic framework for subopts and mfe ($RT \rightarrow \infty$)
- ⇒ Very easy to embed into any existing DP equations

Partition function: Beyond m.f.e. hypothesis

Partition function/Boltzmann probability

- Let ω be an RNA sequence
- \mathcal{S}_ω be the set of sequences compatible with ω ,

$$\text{Partition function } Z_\omega = \sum_{S \in \mathcal{S}_\omega} e^{\frac{-E_{S,\omega}}{RT}}$$

where T is temperature in Kelvin and R is the universal gas constant.

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

$P_{S,\omega}$ is the probability of observing ω in conformation S .

- ⇒ Offers a more dynamic view of the folding process
- ⇒ Gives a model for computing various probabilities (BP, Motifs ...)
- ⇒ Unified algorithmic framework for subopts and mfe ($RT \rightarrow \infty$)
- ⇒ Very easy to embed into any existing DP equations

Partition function: Beyond m.f.e. hypothesis

Partition function/Boltzmann probability

- Let ω be an RNA sequence
- \mathcal{S}_ω be the set of sequences compatible with ω ,

$$\text{Partition function } Z_\omega = \sum_{S \in \mathcal{S}_\omega} e^{\frac{-E_{S,\omega}}{RT}}$$

where T is temperature in Kelvin and R is the universal gas constant.

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

$P_{S,\omega}$ is the probability of observing ω in conformation S .

- ⇒ Offers a more dynamic view of the folding process
- ⇒ Gives a model for computing various probabilities (BP, Motifs ...)
- ⇒ Unified algorithmic framework for subopts and mfe ($RT \rightarrow \infty$)
- ⇒ Very easy to embed into any existing DP equations

Partition function: Beyond m.f.e. hypothesis

From m.f.e. folding to partition function [McC90]:

- Atomic energy increment $E \rightarrow$ Boltzmann factor $e^{-\frac{E}{RT}}$
- Energies contr. move to the exponent:
Sums (+) \rightarrow Products (\times)
- Summing instead of minimizing: Min \rightarrow Sums (Σ)

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

Message #2

From **unambiguous description** partition function (and then statistical sampling) is just one algebra switch (Min, +) \rightarrow (+, \times) away.

Partition function: Beyond m.f.e. hypothesis

From m.f.e. folding to partition function [McC90]:

- Atomic energy increment $E \rightarrow$ Boltzmann factor $e^{\frac{-E}{RT}}$
- Energies contr. move to the exponent:
Sums (+) \rightarrow Products (\times)
- Summing instead of minimizing: Min \rightarrow Sums (\sum)

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

Message #2

From **unambiguous description** partition function (and then statistical sampling) is just one algebra switch (Min, +) \rightarrow (+, \times) away.

Partition function: Beyond m.f.e. hypothesis

From m.f.e. folding to partition function [McC90]:

- Atomic energy increment $E \rightarrow$ Boltzmann factor $e^{\frac{-E}{RT}}$
- Energies contr. move to the exponent:
Sums (+) \rightarrow Products (\times)
- Summing instead of minimizing: Min \rightarrow Sums (\sum)

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

Message #2

From **unambiguous description** partition function (and then statistical sampling) is just one algebra switch (Min, +) \rightarrow (+, \times) away.

Partition function: Beyond m.f.e. hypothesis

From m.f.e. folding to partition function [McC90]:

- Atomic energy increment $E \rightarrow$ Boltzmann factor $e^{\frac{-E}{RT}}$
- Energies contr. move to the exponent:
Sums (+) \rightarrow Products (\times)
- Summing instead of minimizing: Min \rightarrow Sums (\sum)

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

Message #2

From **unambiguous description** partition function (and then statistical sampling) is just one algebra switch (Min, +) \rightarrow (+, \times) away.

Partition function: Beyond m.f.e. hypothesis

From m.f.e. folding to partition function [McC90]:

- Atomic energy increment $E \rightarrow$ Boltzmann factor $e^{\frac{-E}{RT}}$
- Energies contr. move to the exponent:
Sums (+) \rightarrow Products (\times)
- Summing instead of minimizing: Min \rightarrow Sums (\sum)

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

$$Z(i, j) = \sum \left(Z(i, k-1) + e^{\frac{-b(k-1)}{RT}} \right) Z^1(k, j)$$

$$Z^1(i, j) = e^{\frac{-b}{RT}} Z^1(i, j-1) + e^{\frac{-c}{RT}} Z'(i, j)$$

Message #2

From **unambiguous description** partition function (and then statistical sampling) is just one algebra switch (Min, +) \rightarrow (+, \times) away.

Statistical sampling through stochastic traceback

Algorithm SFold [DL03]:

- 1 Generate a random number in $[0, Z'(i, j))$
- 2 Subtract to r individual contributions to $Z'(i, j)$, until $r < 0$
- 3 Recurse over substructures

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

Statistical sampling through stochastic traceback

Algorithm SFold [DL03]:

- 1 Generate a random number in $[0, Z'(i, j))$
- 2 Subtract to r individual contributions to $Z'(i, j)$, until $r < 0$
- 3 Recurse over substructures

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

Statistical sampling through stochastic traceback

Algorithm SFold [DL03]:

- 1 Generate a random number in $[0, Z'(i, j))$
- 2 Subtract to r individual contributions to $Z'(i, j)$, until $r < 0$
- 3 Recurse over substructures

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

\boxed{r}
↓

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

After $\Theta(n)$ operations, recurse over size $n-1$ interval
 \Rightarrow Worst-case time complexity for k samples in $\mathcal{O}(n^2 k)$

Remark: This is a weighted instance of the so-called recursive random generation of decomposable objects [DRT00].

Statistical sampling through stochastic traceback

Algorithm SFold [DL03]:

- 1 Generate a random number in $[0, Z'(i, j))$
- 2 Subtract to r individual contributions to $Z'(i, j)$, until $r < 0$
- 3 Recurse over substructures

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

\boxed{r}
↓

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

After $\Theta(n)$ operations, recurse over size $n - 1$ interval
 \Rightarrow Worst-case time complexity for k samples in $\mathcal{O}(n^2 k)$

Remark: This is a weighted instance of the so-called recursive random generation of decomposable objects [DRT00].

Statistical sampling through stochastic traceback

Algorithm SFold [DL03]:

- 1 Generate a random number in $[0, Z'(i, j))$
- 2 Subtract to r individual contributions to $Z'(i, j)$, until $r < 0$
- 3 Recurse over substructures

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

\boxed{r}
↓

A₁ | A₂ | B_i | B_{i+1} | ... | B_{j-1} | B_j | C_i | C_{i+1} | ... | C_{j-1} | C_j

After $\Theta(n)$ operations, recurse over size $n - 1$ interval
 \Rightarrow Worst-case time complexity for k samples in $\mathcal{O}(n^2 k)$

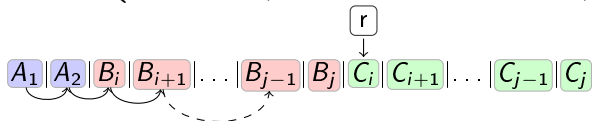
Remark: This is a weighted instance of the so-called recursive random generation of decomposable objects [DRT00].

Statistical sampling through stochastic traceback

Algorithm SFold [DL03]:

- 1 Generate a random number in $[0, \mathcal{Z}'(i, j))$
- 2 Subtract to r individual contributions to $\mathcal{Z}'(i, j)$, until $r < 0$
- 3 Recurse over substructures

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



After $\Theta(n)$ operations, recurse over size $n-1$ interval
 \Rightarrow Worst-case time complexity for k samples in $\mathcal{O}(n^2k)$

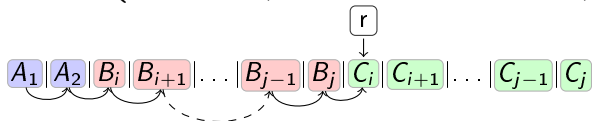
Remark: This is a weighted instance of the so-called recursive random generation of decomposable objects [DRT00].

Statistical sampling through stochastic traceback

Algorithm SFold [DL03]:

- 1 Generate a random number in $[0, \mathcal{Z}'(i, j))$
- 2 Subtract to r individual contributions to $\mathcal{Z}'(i, j)$, until $r < 0$
- 3 Recurse over substructures

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



After $\Theta(n)$ operations, recurse over size $n-1$ interval
 \Rightarrow Worst-case time complexity for k samples in $\mathcal{O}(n^2 k)$

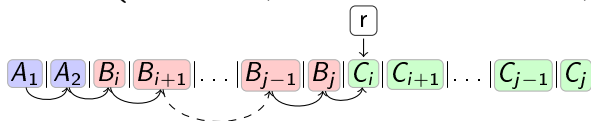
Remark: This is a weighted instance of the so-called recursive random generation of decomposable objects [DRT00].

Statistical sampling through stochastic traceback

Algorithm SFold [DL03]:

- 1 Generate a random number in $[0, \mathcal{Z}'(i, j))$
- 2 Subtract to r individual contributions to $\mathcal{Z}'(i, j)$, until $r < 0$
- 3 Recurse over substructures

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



After $\Theta(n)$ operations, recurse over size $n - 1$ interval
 \Rightarrow Worst-case time complexity for k samples in $\mathcal{O}(n^2 k)$

Remark: This is a weighted instance of the so-called recursive random generation of decomposable objects [DRT00].

Efficient statistical sampling

How to improve statistical sampling?

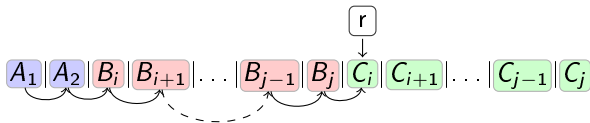
- Improve time complexity:

Average-case time complexity in $\Theta(kn\sqrt{n})$ [Pon08]

($\Theta(n^2)$ arises from recursing on $n - \mathcal{O}(1)$ after $\Theta(n)$ ops)

- Interleaving Bulges (B) and Multiloops (C) contributions
- Boustrophedon [FZV94]
Investigate uneven decompositions first, then even ones !

- Non-redundant generation



Message #3

Boustrophedon search saves $\Theta(\frac{n}{\log n})/\Omega(\frac{\sqrt{n}}{\log n})$ worst/average case.

Efficient statistical sampling

How to improve statistical sampling?

- Improve time complexity:

Average-case time complexity in $\Theta(kn\sqrt{n})$ [Pon08]

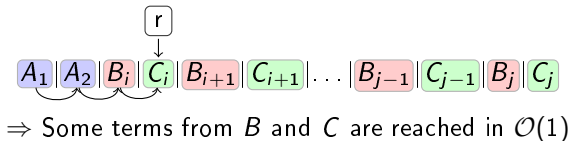
($\Theta(n^2)$ arises from recursing on $n - \mathcal{O}(1)$ after $\Theta(n)$ ops)

- Interleaving Bulges (B) and Multiloops (C) contributions

- Boustrophedon [FZV94]

Investigate uneven decompositions first, then even ones !

- Non-redundant generation



Message #3

Boustrophedon search saves $\Theta(\frac{n}{\log n})/\Omega(\frac{\sqrt{n}}{\log n})$ worst/average case.

Efficient statistical sampling

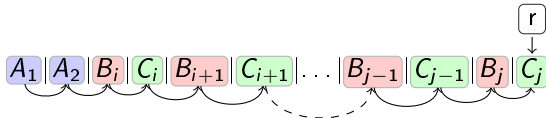
How to improve statistical sampling?

- Improve time complexity:

Average-case time complexity in $\Theta(kn\sqrt{n})$ [Pon08]

($\Theta(n^2)$ arises from recursing on $n - \mathcal{O}(1)$ after $\Theta(n)$ ops)

- Interleaving Bulges (B) and Multiloops (C) contributions
- Boustrophedon [FZV94]
Investigate uneven decompositions first, then even ones !
- Non-redundant generation



\Rightarrow Some terms from B and C are reached in $\mathcal{O}(1)$

But still $\Theta(n^2)$, since $\mathcal{Z}'(i, j) \rightarrow (\mathcal{Z}'(i + 1, k - 1), \mathcal{Z}^1(k, j - 1))$

Message #3

Boustrophedon search saves $\Theta(\frac{n}{\log n})/\Omega(\frac{\sqrt{n}}{\log n})$ worst/average case.

Efficient statistical sampling

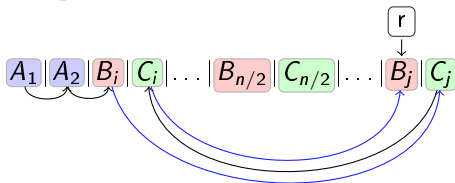
How to improve statistical sampling?

- Improve time complexity:

Average-case time complexity in $\Theta(kn\sqrt{n})$ [Pon08]

($\Theta(n^2)$ arises from recursing on $n - \mathcal{O}(1)$ after $\Theta(n)$ ops)

- Interleaving Bulges (B) and Multiloops (C) contributions
 - **Boustrophedon** [FZV94]
 - Investigate uneven decompositions first, then even ones !
- Non-redundant generation



Message #3

Boustrophedon search saves $\Theta(\frac{n}{\log n})/\Omega(\frac{\sqrt{n}}{\log n})$ worst/average case.

Efficient statistical sampling

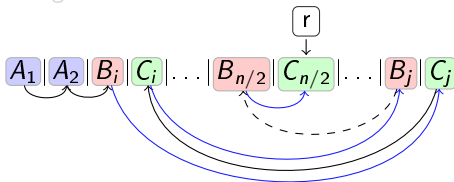
How to improve statistical sampling?

- Improve time complexity:

Average-case time complexity in $\Theta(kn\sqrt{n})$ [Pon08]

($\Theta(n^2)$ arises from recursing on $n - \mathcal{O}(1)$ after $\Theta(n)$ ops)

- Interleaving Bulges (B) and Multiloops (C) contributions
- Boustrophedon [FZV94] $\Rightarrow \Theta(n \log(n))$ worst-case
Investigate uneven decompositions first, then even ones !
- Non-redundant generation



Worst-case: Divide evenly at each step [GK81] $\Rightarrow \Theta(n \log(n))$

Message #3

Boustrophedon search saves $\Theta(\frac{n}{\log n}) / \Omega(\frac{\sqrt{n}}{\log n})$ worst/average case.

Efficient statistical sampling

How to improve statistical sampling?

- Improve time complexity:

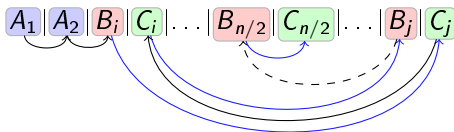
Average-case time complexity in $\Theta(kn\sqrt{n})$ [Pon08]

($\Theta(n^2)$ arises from recursing on $n - \mathcal{O}(1)$ after $\Theta(n)$ ops)

- Interleaving Bulges (B) and Multiloops (C) contributions
- Boustrophedon [FZV94]

Investigate uneven decompositions first, then even ones !

- Non-redundant generation



Worst-case: Divide evenly at each step [GK81] $\Rightarrow \Theta(n \log(n))$

Message #3

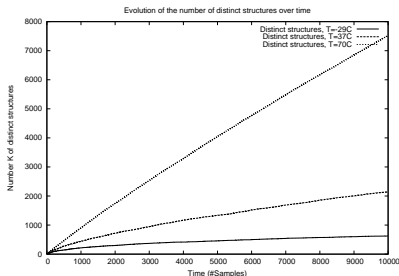
Boustrophedon search saves $\Theta(\frac{n}{\log n})/\Omega(\frac{\sqrt{n}}{\log n})$ worst/average case.

Efficient statistical sampling

How to improve statistical sampling?

- Improve time complexity
- Non-redundant generation (with **D. Gardy** and A. Lorenz)

For each sampled structure one can compute the actual probability
⇒ It does not make any sense to sample it twice!



How long will it take to get k **distinct** samples?

Efficient statistical sampling

How to improve statistical sampling?

- Improve time complexity
- Non-redundant generation (with **D. Gardy** and A. Lorenz)

How long will it take to get k distinct samples?

Full collection ($k = \#structures$): $E[C] \approx \mathcal{Z}' \cdot n$

Way larger than $\#structures \Rightarrow$ Exponential number of collisions.

Numerical values (Homopolymer/Nussinov energy/ $T=37$):

$$E[C] \sim K \cdot 4.332^n / \sqrt{n} \text{ and } \#structures: S_n \sim K' \cdot 2.618^n / n\sqrt{n}$$

\Rightarrow Each structure is sampled $1.65^n \cdot n$ times ($\neq \Theta(n)$ uniform dist.)

Message #4

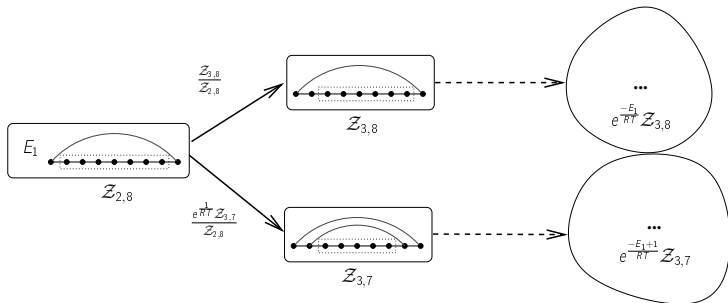
For any RNA there exists k such that the time for sampling k **distinct** sec. str. is **heavily dominated** by the cost of collisions.

k depends on the length \Rightarrow Still need to push further our analysis...

Efficient statistical sampling

How to improve statistical sampling?

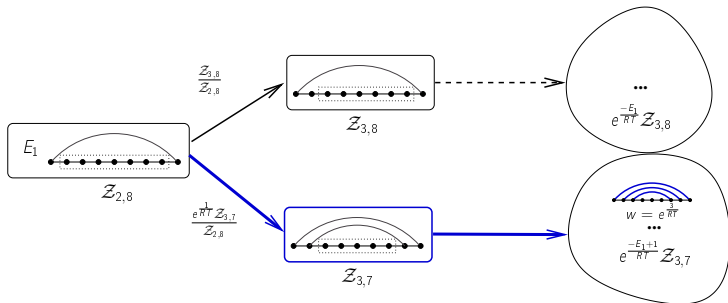
- Improve time complexity
- Non-redundant generation (with D. Gardy and A. Lorenz)
 - Build **prefix tree** for parse traces, storing in each node the contributions $K = \sum_{S \in \mathcal{R}} e^{-\frac{E_S}{RT}}$ of already sampled structures \mathcal{R}
 - During traceback, **modify contributions** of terms using K [Pon08]



Efficient statistical sampling

How to improve statistical sampling?

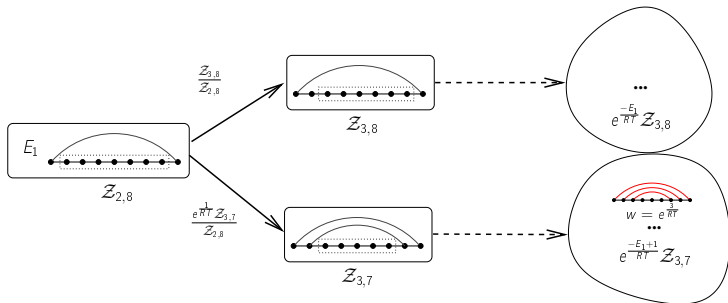
- Improve time complexity
- Non-redundant generation (with D. Gardy and A. Lorenz)
 - Build **prefix tree** for parse traces, storing in each node the contributions $K = \sum_{S \in \mathcal{R}} e^{-\frac{E_S}{RT}}$ of already sampled structures \mathcal{R}
 - During traceback, **modify contributions** of terms using K [Pon08]



Efficient statistical sampling

How to improve statistical sampling?

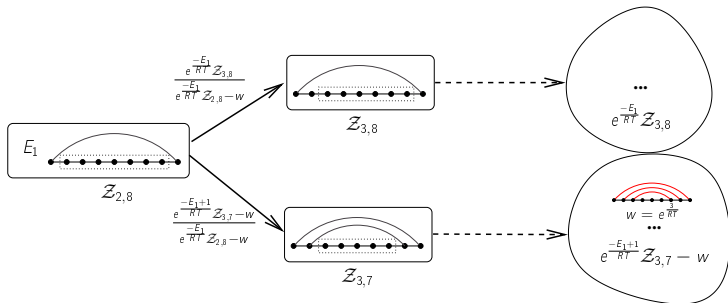
- Improve time complexity
- Non-redundant generation (with D. Gardy and A. Lorenz)
 - Build **prefix tree** for parse traces, storing in each node the contributions $K = \sum_{S \in \mathcal{R}} e^{-\frac{E_S}{RT}}$ of already sampled structures \mathcal{R}
 - During traceback, **modify contributions** of terms using K [Pon08]



Efficient statistical sampling

How to improve statistical sampling?

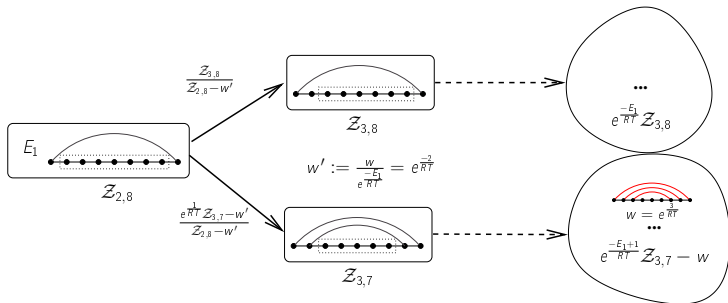
- Improve time complexity
- Non-redundant generation (with D. Gardy and A. Lorenz)
 - Build **prefix tree** for parse traces, storing in each node the contributions $K = \sum_{S \in \mathcal{R}} e^{-\frac{E_S}{RT}}$ of already sampled structures \mathcal{R}
 - During traceback, **modify contributions** of terms using K [Pon08]



Efficient statistical sampling

How to improve statistical sampling?

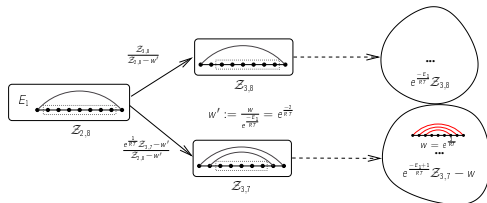
- Improve time complexity
- Non-redundant generation (with D. Gardy and A. Lorenz)
 - Build **prefix tree** for parse traces, storing in each node the contributions $K = \sum_{S \in \mathcal{R}} e^{-\frac{E_S}{RT}}$ of already sampled structures \mathcal{R}
 - During traceback, **modify contributions** of terms using K [Pon08]



Efficient statistical sampling

How to improve statistical sampling?

- Improve time complexity
- Non-redundant generation (with D. Gardy and A. Lorenz)
 - Build **prefix tree** for parse traces, storing in each node the contributions $K = \sum_{S \in \mathcal{R}} e^{\frac{-E_S}{RT}}$ of already sampled structures \mathcal{R}
 - During traceback, **modify contributions** of terms using K [Pon08]



Message #5

Storing parse trees and biasing local choices, one can perform non-redundant sampling in $\mathcal{O}(kn \log(n))$ time.

- Combinatorics gives a convenient framework for validating/analyzing/improving dynamic programming algorithms
- Statistical sampling = (Weighted) random generation of combinatorial structures + constraints
- During stochastic traceback, reordering comparisons saves time!
- One does not benefit from redundancy \Rightarrow Non-redundant sampling

Open questions:

- When do collisions overcome the complexity of sampling?
- Does there exist sequential alternatives to RNASubopt?

Thanks for E. Rivas and E. Westhof!!!



Y. Ding and E. Lawrence.

A statistical sampling algorithm for RNA secondary structure prediction.
Nucleic Acids Research, 31(24):7280–7301, 2003.



A. Denise, O. Roques, and M. Termier.

Random generation of words of context-free languages according to the frequencies of letters.
In D. Gardy and A. Mokedem, editors, *Mathematics and Computer Science: Algorithms, Trees, Combinatorics and probabilities*, Trends in Mathematics, pages 113–125. Birkhäuser, 2000.



P. Flajolet, P. Zimmermann, and B. Van Cutsem.

Calculus for the random generation of labelled combinatorial structures.
Theoretical Computer Science, 132:1–35, 1994.



D. H. Greene and D. E. Knuth.

Mathematics for the Analysis of Algorithms.
Birkhauser Boston, 1981.



J.S. McCaskill.

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



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.



M. S. Waterman.

Secondary structure of single stranded nucleic acids.
Advances in Mathematics Supplementary Studies, 1(1):167–212, 1978.