

Combinatorial insight into RNA structure

Analysis and Design of *ab-initio* RNA algorithms

Yann Ponty

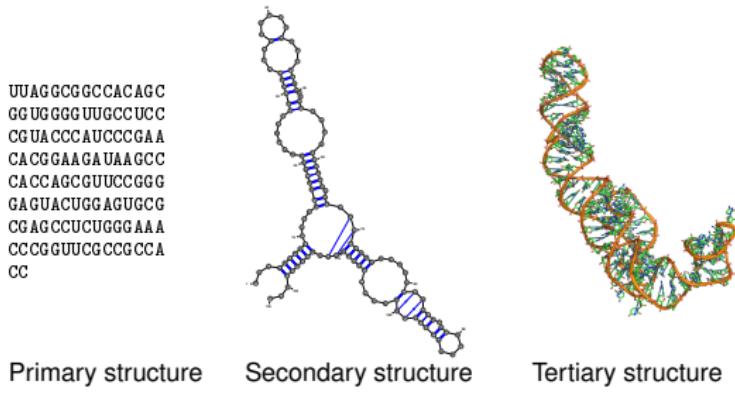
Bioinformatics Team
École Polytechnique/CNRS/INRIA AMIB – France

April 24, 2011

- Part 1: Foreword/Generating functions 101
- Part 2: Enumeration of RNA shapes
- Part 3: Realistic models for random RNA structures
- Part 4: Improved statistical sampling
- Part 5: Conclusion

Part I

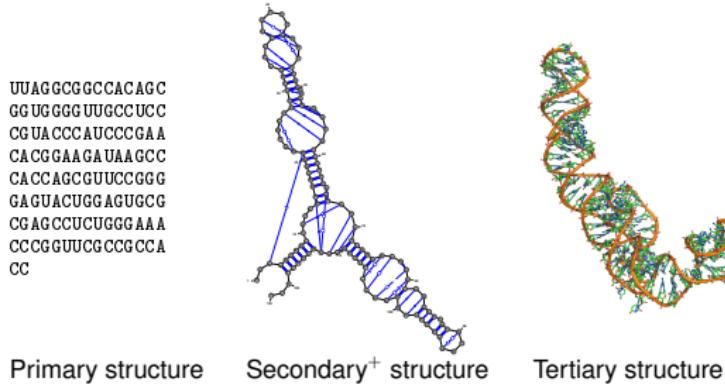
Foreword/Generating functions 101



Source: 5s rRNA (PDBID: 1K73-B)

Definition

Secondary structures of RNA =
Maximal non-crossing subset of canonical base-pairs.



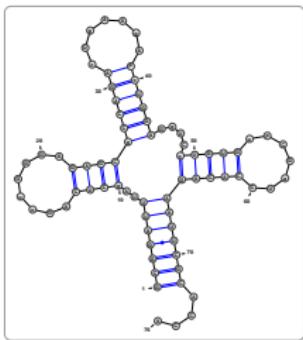
Source: 5s rRNA (PDBID: 1K73-B)

Definition

Secondary structures of RNA =

Maximal non-crossing subset of **canonical** base-pairs.

Various representations for a versatile molecule



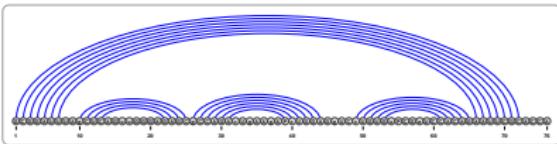
Outer planar graph

(((((.,((.,.....)))) (((((.....))))))...(((.,.....))))....)

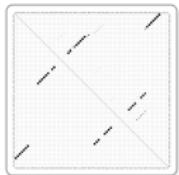
Well-parenthesized expression



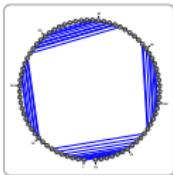
Mountain view



Non intersecting arcs



Dot plot



Feynman's diagram

Different objects
yet
Common combinatorial structure

Minimum Free-Energy (MFE)

Functional folding of an RNA = Minimum free-energy structure.

Nussinov/Jacobson (NJ) energy model: Simplest *nearest-neighbor* model.

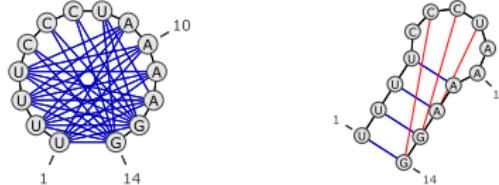
$$\text{Free-energy} = -\#\text{Base-pairs}$$

Folding in NJ \Leftrightarrow Maximizing number of base-pairs.

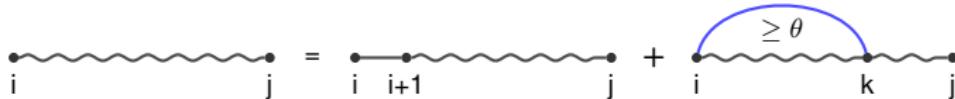
(\Leftrightarrow Max independent set in circle graphs \Leftrightarrow weighted CYK ...)

Example:

UUUUUCCCUAAAAGG



Alt. : Assign weight to each base-pairs based on number of hydrogen bonds
 $\Delta G(G \equiv C) = -3$ $\Delta G(A = U) = -2$ $\Delta G(G - U) = -1$



Recurrence on the Minimal Free-Energy $N_{i,j}$ of any subinterval of an RNA :

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

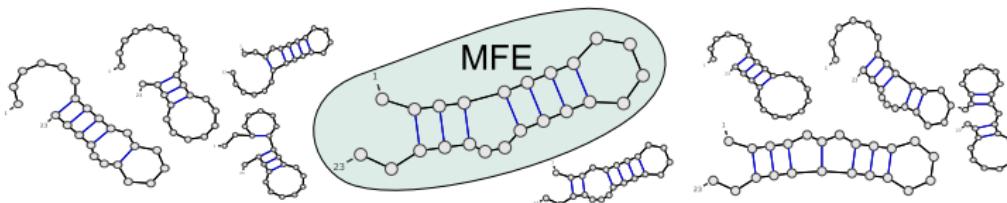
$$N_{i,j} = \min \left\{ \begin{array}{ll} N_{i+1,j} & (i \text{ unpaired}) \\ \min_{k=i+\theta+1}^j E_{i,j} + N_{i+1,k-1} + N_{k+1,j} & (i \text{ comp. with } k) \end{array} \right.$$

After computing N , the choices elected by \min allow to rebuild (backtrack over) the MFE structure.

⇒ Dynamic programming for folding [NJ80].

Remark: Overly simplistic, but more complex energy models will typically already predict 73% of experimentally-determined base-pairs...

RNA *breathes* \Rightarrow There is not necessarily one functional conformation!



Boltzmann-ensemble paradigm

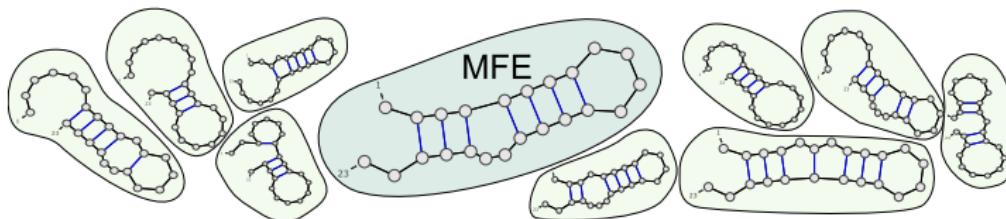
All possible secondary structures for an RNA coexist in a probability distribution, the Boltzmann distribution.

Bonus: MFE structure probability can be significantly remote, pointing toward unstructured RNAs.

Conversely, base-pairs having probability $> .99$ are valid 95% of the time.

- \Rightarrow Similar structures can bundle up.
- \Rightarrow Functional structures must be sought in sub-optimal structures.

RNA *breathes* ⇒ There is not necessarily one functional conformation!



Boltzmann-ensemble paradigm

All possible secondary structures for an RNA coexist in a probability distribution, the Boltzmann distribution.

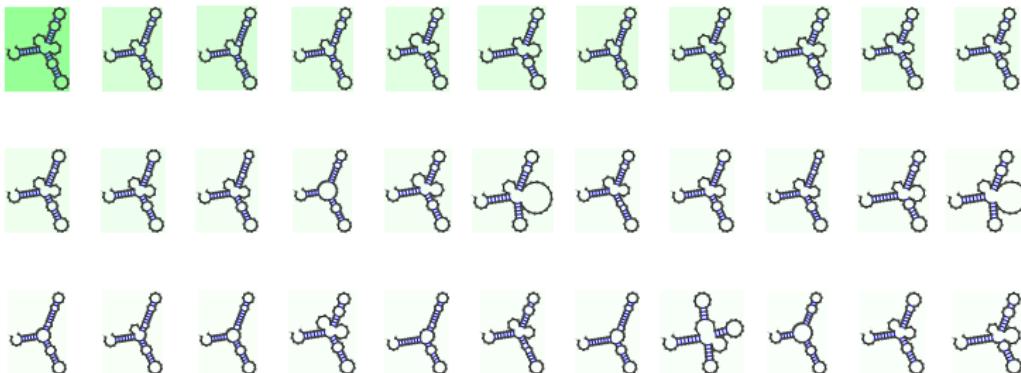
Bonus: MFE structure probability can be significantly remote, pointing toward unstructured RNAs.

Conversely, base-pairs having probability $> .99$ are valid 95% of the time.

- ⇒ Similar structures can bundle up.
- ⇒ Functional structures must be sought in sub-optimal structures.

Why use combinatorics?

Boltzmann ensemble is a (weighted) combinatorial class.



Studying it as such *cleans out the details* and helps:

- Access **asymptotic properties** of secondary structures
- Investigate worst and average-case complexities
- Design **better algorithms and methods** for RNA

Let $| \cdot |$ be a size function over objects (Sequences, trees, ...).

Combinatorial classes are (infinite) sets \mathcal{C} of objects whose restrictions \mathcal{C}_n to objects of size n are of finite cardinality.

Definition (Generating functions)

Let \mathcal{C} be a combinatorial class and $c_n = |\mathcal{C}_n|$ the number of objects of size n in \mathcal{C} , then the generating function of \mathcal{C} is $C(z)$ s. t.

$$C(z) = \sum_{s \in \mathcal{C}} z^{|s|} = \sum_{n \geq 0} c_n z^n$$

Closed forms for $C(z)$ are often easy to find ...

DNA example: $\mathcal{D} := \{a, c, g, t\}^*$ $\Rightarrow d_n = 4^n$

$$\text{and } C(z) = 1 + 4z + 16z^2 + 64z^3 + \dots = \sum_{n \geq 0} 4^n z^n = \frac{1}{1-4z}$$

... and very often much simpler than for c_n !!!

From a **class specification**, one can directly establish the gen. fun.

Historically on languages, from Schützenberger's observation that

Gen. fun. are commutative images of languages

Grammar	Generating function	Coefficients
$C \rightarrow \varepsilon$	$C(z) = z^0 = 1$	$c_n = \mathbb{1}_{\{0\}}(n)$
$C \rightarrow t$	$C(z) = z^1 = z$	$c_n = \mathbb{1}_{\{1\}}(n)$
$C \rightarrow A \mid B$	$C(z) = A(z) + B(z)$	$c_n = a_n + b_n$
$C \rightarrow A.B$	$C(z) = A(z) \cdot B(z)$	$c_n = \sum_{k=0}^n a_k b_{n-k}$

Remark: One needs to ensure that unions are disjoint and concatenations unambiguous.

DNA example : $\{a, c, g, t\}^* \Leftrightarrow D \rightarrow a.D \mid c.D \mid g.D \mid t.D \mid \varepsilon$

$$\Rightarrow D(z) = z \cdot D(z) + z \cdot D(z) + z \cdot D(z) + z \cdot D(z) + 1$$

From a **class specification**, one can directly establish the gen. fun.

Historically on languages, from Schützenberger's observation that

Gen. fun. are commutative images of languages

Grammar	Generating function	Coefficients
$C \rightarrow \varepsilon$	$C(z) = z^0 = 1$	$c_n = \mathbb{1}_{\{0\}}(n)$
$C \rightarrow t$	$C(z) = z^1 = z$	$c_n = \mathbb{1}_{\{1\}}(n)$
$C \rightarrow A \mid B$	$C(z) = A(z) + B(z)$	$c_n = a_n + b_n$
$C \rightarrow A.B$	$C(z) = A(z) \cdot B(z)$	$c_n = \sum_{k=0}^n a_k b_{n-k}$

Remark: One needs to ensure that unions are disjoint and concatenations unambiguous.

DNA example : $\{a, c, g, t\}^* \Leftrightarrow D \rightarrow a.D \mid c.D \mid g.D \mid t.D \mid \varepsilon$

$$\Rightarrow D(z) = 4z \cdot D(z) + 1$$

From a **class specification**, one can directly establish the gen. fun.

Historically on languages, from Schützenberger's observation that

Gen. fun. are commutative images of languages

Grammar	Generating function	Coefficients
$C \rightarrow \varepsilon$	$C(z) = z^0 = 1$	$c_n = \mathbb{1}_{\{0\}}(n)$
$C \rightarrow t$	$C(z) = z^1 = z$	$c_n = \mathbb{1}_{\{1\}}(n)$
$C \rightarrow A \mid B$	$C(z) = A(z) + B(z)$	$c_n = a_n + b_n$
$C \rightarrow A.B$	$C(z) = A(z) \cdot B(z)$	$c_n = \sum_{k=0}^n a_k b_{n-k}$

Remark: One needs to ensure that unions are disjoint and concatenations unambiguous.

DNA example : $\{a, c, g, t\}^* \Leftrightarrow D \rightarrow a.D \mid c.D \mid g.D \mid t.D \mid \varepsilon$

$$\Rightarrow D(z) = \frac{1}{1 - 4z}$$

Disclaimer

What follows, although true in this context, is embarrassingly simplistic.
A rigorous presentation can (and must) be found in Flaj./Sedg. 08.

A **singularity** is a point $z = \rho$ where $C(z)$ is no longer analytic.
Asymptotics of coeff c_n are driven by the **singularities** of $C(z)$.

1st principle

Location of the dominant (smallest) singularity ρ dictates the exponential growth $\Rightarrow \frac{c_n}{\rho^{-n}} = o(\alpha^n)$, $\forall \alpha > 1$.

DNA example: $D(z) = 1/(1 - 4z) \Rightarrow \rho = 1/4 \Rightarrow d_n \sim 4^n P(n)$.

2nd principle

Nature of ρ dictates subexponential part $P(n)$ s.t. $c_n \sim \rho^{-n} P(n)$.

Basic scale: If one can rewrite $C(z)$ as

$$C(z) = f(z) + g(z)(1 - z/\rho)^\alpha$$

where f and g are analytic $\forall |z| < |\rho|$ and non-null at ρ , then

$$c_n \equiv [z^n] C(z) \sim \frac{g(\rho)\rho^{-n}}{\Gamma(-\alpha)n^{\alpha+1}}$$

Example: $D(z) = \frac{1}{1-4z} \Rightarrow c_n \sim 4^n$
 $(\rho = 1/4, \alpha = -1, f(z) = 0, \text{ and } g(z) = 1)$

Asymptotic estimates are obtained using a 4 steps meta-algorithm:

1

Find the right model

2

Translate into grammar

3

Translate into system and solve g. f.

4

Singularity analysis yields asymptotics

Asymptotic estimates are obtained using a 4 steps meta-algorithm:

1

Find the right model

2

Translate into grammar

3

Translate into system and solve g. f.

4

Singularity analysis yields asymptotics

Asymptotic estimates are obtained using a 4 steps meta-algorithm:

1

Find the right model

2

Translate into grammar

3

Translate into system and solve g. f.

4

Singularity analysis yields asymptotics

Asymptotic estimates are obtained using a 4 steps meta-algorithm:

- 1 Find the right model
- 2 Translate into grammar
- 3 Translate into system and solve g. f.
- 4 Singularity analysis yields asymptotics

Appetizer: Motzkin words

Let us count dot-bracket notations (Motzkin words)

1

$$\bullet \text{---} \text{---} \text{---} \bullet = \bullet \cdot \bullet \text{---} \text{---} \text{---} \bullet \vee \bullet \text{---} \text{---} \text{---} \bullet \text{---} \text{---} \bullet \vee \varepsilon$$

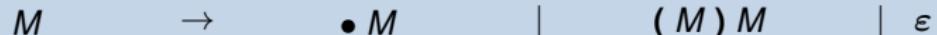
Appetizer: Motzkin words

Let us count dot-bracket notations (Motzkin words)

1



2



Appetizer: Motzkin words

Let us count dot-bracket notations (Motzkin words)

$$1 \quad \bullet \text{---} \bullet = \bullet \cdot \bullet \text{---} \bullet \vee \bullet \text{---} \bullet \text{---} \bullet \vee \varepsilon$$

$$2 \quad M \rightarrow \bullet M \quad | \quad (M)M \quad | \quad \varepsilon$$

$$3 \quad M(z) = z \cdot M(z) + z \cdot M(z) \cdot z \cdot M(z) + 1$$

$$= \frac{1-z \pm \sqrt{1-2z-3z^2}}{2z^2}$$

Appetizer: Motzkin words

Let us count dot-bracket notations (Motzkin words)

1  =   $\vee \varepsilon$

2 M \rightarrow $\bullet M$ $|$ $(M)M$ $|$ ε

3 $M(z) = z \cdot M(z) + z \cdot M(z) \cdot z \cdot M(z) + 1$

$$= \begin{cases} \frac{1-z+\sqrt{1-2z-3z^2}}{2z^2} = \frac{1}{z^2} - \frac{1}{z} - 1 - z - 2z^2 + \mathcal{O}(z^3) \\ \frac{1-z-\sqrt{1-2z-3z^2}}{2z^2} = 1 + z + 2z^2 + 4z^3 + 9z^4 + \mathcal{O}(z^5) \end{cases}$$

Appetizer: Motzkin words

Let us count dot-bracket notations (Motzkin words)

1  =  \vee  $\vee \varepsilon$

2 M \rightarrow $\bullet M$ $|$ $(M)M$ $|$ ε

3 $M(z) = z \cdot M(z) + z \cdot M(z) \cdot z \cdot M(z) + 1$

$$= \frac{1-z-\sqrt{1-2z-3z^2}}{2z^2}$$

4 $\rho = 1/3$, $M(z) = \frac{1-z}{2z^2} - g(z) \cdot \sqrt{1-z/\rho}$, and $g(z) := \frac{\sqrt{1+z}}{2z^2}$

$$\Rightarrow s_n \equiv [z^n]M(z) \sim \frac{g(\rho)\rho^{-n}}{\Gamma(-\alpha)n^{\alpha+1}} = \frac{3\sqrt{3}}{2\sqrt{\pi}} \cdot \frac{3^n}{n\sqrt{n}}(1 + \mathcal{O}(1/n))$$

RNA secondary structures

Consider RNA secondary structures (Waterman 78)

1



RNA secondary structures

Consider RNA secondary structures (Waterman 78)

1



2

$$S \rightarrow \bullet S | (S_{>0}) S | \epsilon$$

RNA secondary structures

Consider RNA secondary structures (Waterman 78)

1 

2 $S \rightarrow \bullet S | (T) S | \epsilon$

$T \rightarrow \bullet S | (T) S$

RNA secondary structures

Consider RNA secondary structures (Waterman 78)

1  ≥ 1

2 $S \rightarrow \bullet S | (T) S | \epsilon$
 $T \rightarrow \bullet S | (T) S$

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

RNA secondary structures

Consider RNA secondary structures (Waterman 78)

1



2

$$S \rightarrow \bullet S | (T) S | \varepsilon$$

$$T \rightarrow \bullet S | (T) S$$

3

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

4

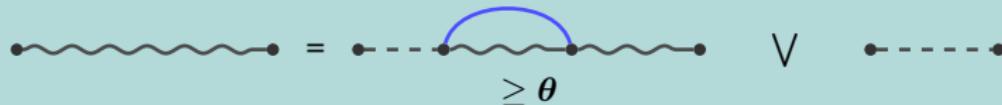
$$\rho = \frac{3-\sqrt{5}}{2} = 1 - \phi$$

$$[z^n]S(z) = \sqrt{\frac{15+7\sqrt{5}}{8\pi}} \cdot \frac{\left(\frac{3+\sqrt{5}}{2}\right)^n}{n\sqrt{n}} (1 + \mathcal{O}(1/n)) \sim 1.1 \cdot \frac{2.6^n}{n\sqrt{n}}$$

RNA secondary structures

Let us generalize the θ constraint

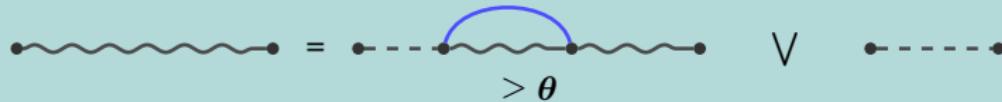
1



RNA secondary structures

Let us generalize the θ constraint

1



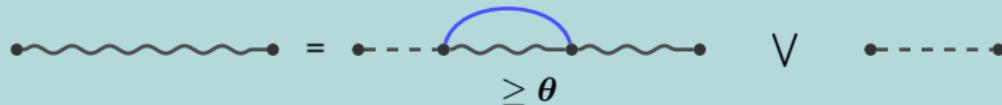
2

$$S \rightarrow U(S_{\geq \theta})S \mid U \quad U \rightarrow \bullet U \mid \epsilon$$

RNA secondary structures

Let us generalize the θ constraint

1



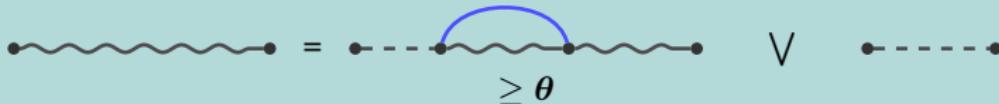
2

$$S \rightarrow U(T) S | U$$
$$T \rightarrow U(T) S | \bullet^\theta U$$
$$U \rightarrow \bullet U | \epsilon$$

RNA secondary structures

Let us generalize the θ constraint

1



2

$$\begin{array}{ll} S \rightarrow U(T)S \mid U & U \rightarrow \bullet U \mid \epsilon \\ T \rightarrow U(T)S \mid \bullet^\theta U & \end{array}$$

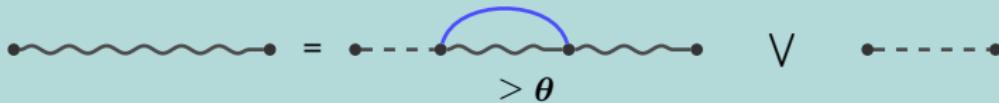
3

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

RNA secondary structures

Let us generalize the θ constraint

1



2

$$\begin{aligned} S &\rightarrow U(T)S \mid U \\ T &\rightarrow U(T)S \mid \bullet^\theta U \end{aligned}$$

3

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

4

$$s_n \sim K \cdot \frac{\beta^n}{n\sqrt{n}} (1 + \mathcal{O}(1/n))$$

θ	0	1	3	10
β	3.	2.62	2.29	2.02

Message #1

Finding the right decomposition (DP) is a combinatorial task.

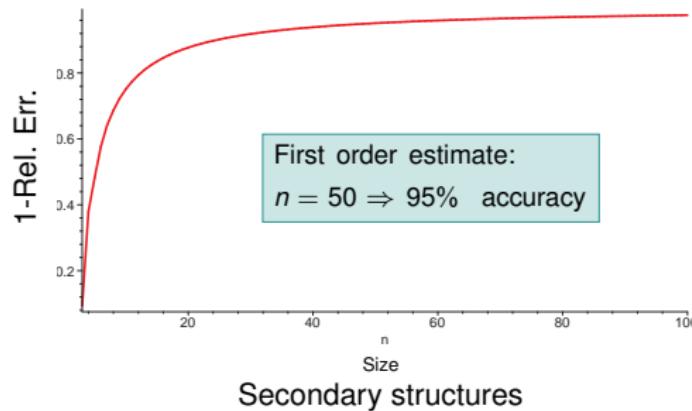
Message #1

Finding the right decomposition (DP) is a combinatorial task.

Message #2

Applying automatic theorems gives precise asymptotic equivalents.

+ Asymptotic regime is reached for small lengths.



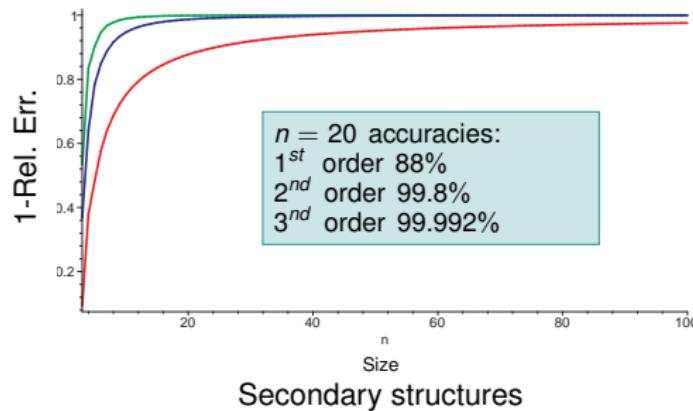
Message #1

Finding the right decomposition (DP) is a combinatorial task.

Message #2

Applying automatic theorems gives precise asymptotic equivalents.

+ Asymptotic regime is reached for small lengths.



Message #1

Finding the right decomposition (DP) is a combinatorial task.

Message #2

Applying automatic theorems gives precise asymptotic equivalents.
+ Asymptotic regime is reached for small lengths.

Message #3

There is a large exponential number of structures of size n :
Homopolymer model: $\Omega(2^n)$ Stickiness model: $\mathcal{O}(1.8^n/n^{3/2})$

Part II

Enumeration of RNA shapes

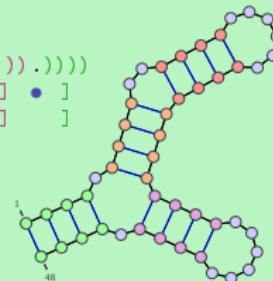
Definition (RNA shapes [Giegerich *et al.*])

Coarse-grain representation hierarchy for RNA sec. struct.

Based on the *underlying backbone* structure.

Example

Sec. str. (((((.(((((((.....))))))))))(((((.....))))..)))
 π' -shape [• [• [[•]]] [[•]] [[•]]]
 π -shape [[- - -]] [[- -]] [[- -]]



Definition (RNA shapes [Giegerich *et al.*])

Coarse-grain representation hierarchy for RNA sec. struct.

Based on the *underlying backbone* structure.

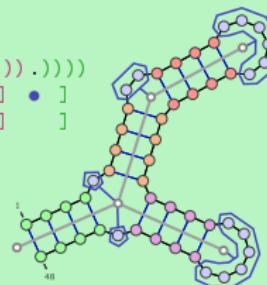
Example

Sec. str. (((((.(((((((.....))))))))))(((((.....))))..))))

π' -shape [• [• [• []]]] [• [] •]]

π -shape [[- - -]] [[]] [[]]]

Contract identical consecutive characters



Definition (RNA shapes [Giegerich *et al.*])

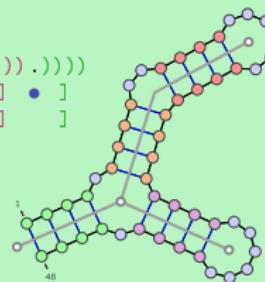
Coarse-grain representation hierarchy for RNA sec. struct.

Based on the *underlying backbone* structure.

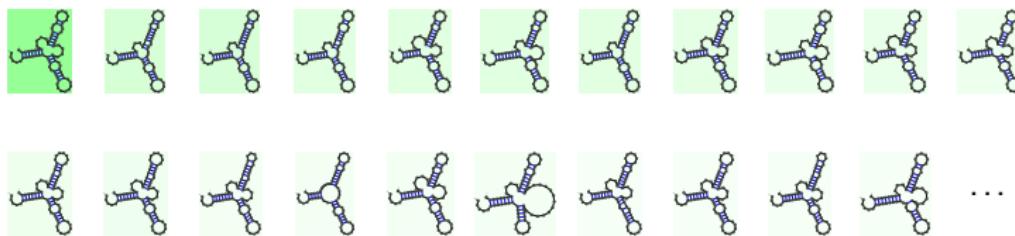
Example

Sec. str. (((((.(((((((.....))))))))))) (((((.....)))).,))))
 π' -shape [• [• [[•]]] [•]] [•]]
 π -shape [[- - -]] [] []]

Remove unpaired regions
Contract nested helices

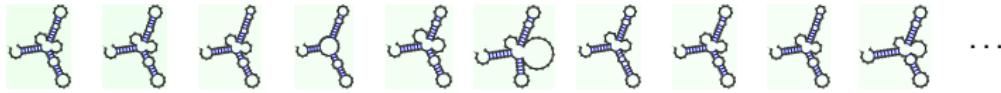
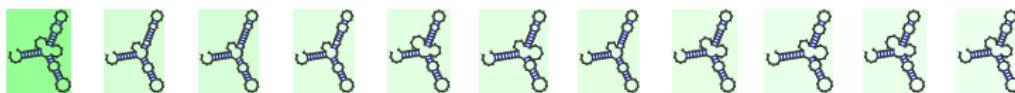


RNA shapes allow a hierarchical search in the Boltzmann ensemble



10000 samples \Rightarrow 1727 Secondary structures...

RNA shapes allow a hierarchical search in the Boltzmann ensemble



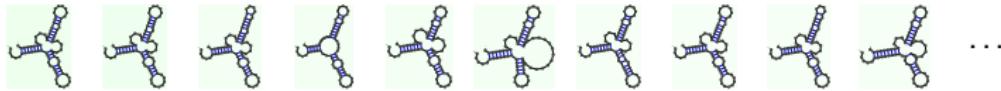
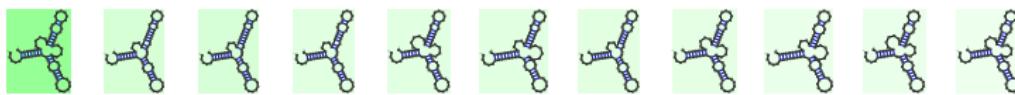
10000 samples \Rightarrow 1727 Secondary structures...



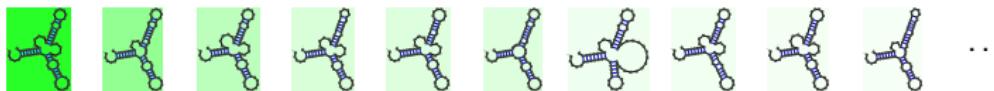
... 406 π' -shapes...

Motivation

RNA shapes allow a hierarchical search in the Boltzmann ensemble



10000 samples \Rightarrow 1727 Secondary structures...

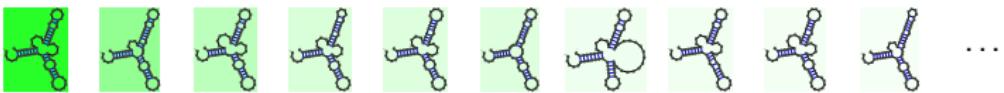
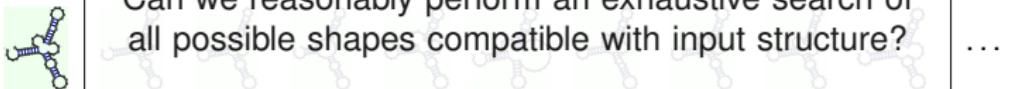
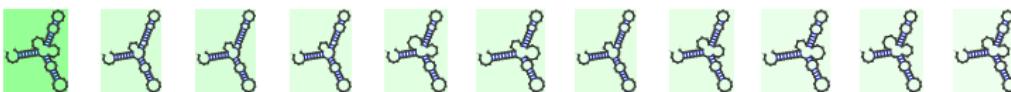


$\dots 406 \pi'$ -shapes...



\dots but only 9 π -shapes!

RNA shapes allow a hierarchical search in the Boltzmann ensemble

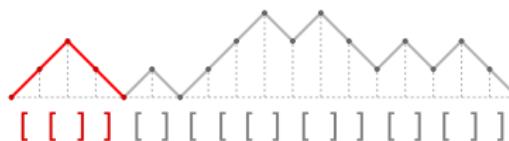


... 406 π' -shapes...



... but only 9 π -shapes!

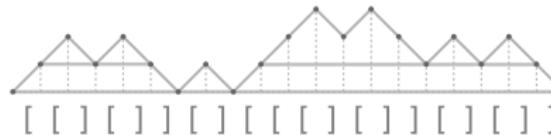
Objective: Count π -shapes with $2n$ parentheses.



1

π -shapes are bracket words avoiding the `[[...]]` motif.

Objective: Count π -shapes with $2n$ parentheses.



1

π -shapes are bracket words avoiding the $[[\dots]]$ motif.

2

$S \rightarrow [S / \{ \dots \}] S \mid [S / \{ \dots \}]$

Objective: Count π -shapes with $2n$ parentheses.



1 π -shapes are bracket words avoiding the $[[\dots]]$ motif.

2 $S \rightarrow [\mathbf{T}]S|[\mathbf{T}]$ $\mathbf{T} \rightarrow [\mathbf{T}]S|\varepsilon$

Objective: Count π -shapes with $2n$ parentheses.

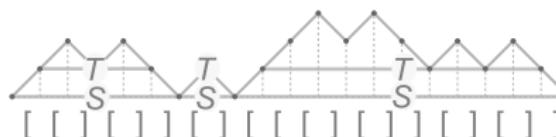


1 π -shapes are bracket words avoiding the $[[\dots]]$ motif.

2 $S \rightarrow [\mathbf{T}]S|[\mathbf{T}]$ $\mathbf{T} \rightarrow [\mathbf{T}]S|\varepsilon$

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

Objective: Count π -shapes with $2n$ parentheses.



1 π -shapes are bracket words avoiding the `[[...]]` motif.

2 $S \rightarrow [T]S|[T]$ $T \rightarrow [T]S|\varepsilon$

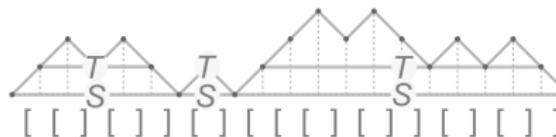
3
$$S(z) = \frac{1 - z^2 - \sqrt{1 - 2z^2 - 3z^4}}{2z^2}$$

4
$$s_{2n} \sim \frac{\sqrt{3}}{2\sqrt{\pi}} \cdot \frac{3^n}{n\sqrt{\pi}} (1 + \mathcal{O}(1/n)) \quad \text{and} \quad s_{2n+1} = 0$$

Remark: Doesn't this look familiar?

These are Motzkin words, i.e. our first example.

Objective: Count π -shapes with $2n$ parentheses.



1 π -shapes are bracket words avoiding the `[[...]]` motif.

2 $S \rightarrow [T]S|[T]$ $T \rightarrow [T]S|\varepsilon$

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

4
$$s_{2n} \sim \frac{\sqrt{3}}{2\sqrt{\pi}} \cdot \frac{3^n}{n\sqrt{\pi}} (1 + \mathcal{O}(1/n)) \quad \text{and} \quad s_{2n+1} = 0$$

Remark: Doesn't this look familiar?

These are Motzkin words, i.e. our first example.

Limitations

Number of π -shapes of size n

\neq

Number of π -shapes compatible with RNA of size n

Reasons:

- 1 Shapes of size $\leq n$ should be considered
- 2 Forming a hairpin loop [] takes at least $\theta + 2$ bases

2 $S \rightarrow [T]S|[T]$ $T \rightarrow [T]S|\varepsilon$

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

4 For n even: $s_{2n} \sim \frac{3\sqrt{3}}{4\sqrt{\pi}} \cdot \frac{3^n}{n\sqrt{\pi}} (1 + \mathcal{O}(1/n)) \approx 0.48 \cdot \frac{3^n}{n\sqrt{n}}$

Limitations

Number of π -shapes of size n

\neq

Number of π -shapes compatible with RNA of size n

Reasons:

- 1 Shapes of size $\leq n$ should be considered
- 2 Forming a hairpin loop [] takes at least $\theta + 2$ bases

2 $S \rightarrow [T]S|[T]$ $T \rightarrow [T]S|\bullet^\theta$

$R \rightarrow \square S | \varepsilon$

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

4
$$r_{2n} \sim \frac{3\sqrt{3}}{4\sqrt{\pi}} \cdot \frac{3^n}{n\sqrt{n}} (1 + \mathcal{O}(1/n)) \Rightarrow r_n \approx 2.07 \cdot \frac{1.73^n}{n\sqrt{n}}$$

Limitations

Number of π -shapes of size n

\neq

Number of π -shapes compatible with RNA of size n

Reasons:

- 1 Shapes of size $\leq n$ should be considered
- 2 Forming a hairpin loop [] takes at least $\theta + 2$ bases

2 $S \rightarrow [T]S|[T]$ $T \rightarrow [T]S|\bullet^\theta$
 $R \rightarrow \square S | \varepsilon$

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

4
$$\theta = 3 \Rightarrow r_n \approx 2.44 \frac{1.32^n}{n\sqrt{n}}$$

Objective: Count π' -shapes compatible with RNA of length n .

1 π' -shapes = bracket words avoiding motifs $[[\dots]]$ and $\bullet\bullet$

2 $R \rightarrow \square R | S \quad S \rightarrow U[T]S | U \quad U \rightarrow \diamond | \varepsilon$

$T \rightarrow U[T]U[T]S | \diamond[T] | [T]\diamond | \diamond[T]\diamond | \bullet^\theta$

3 $\theta = 3, R(z) = \frac{1+2z^2+2z^3+z^4-z^5-z^6-\sqrt{1-4z^3-2z^4-2z^5+2z^6-7z^8-z^{10}+2z^{11}+z^{12}}}{2z^2(1-z^2)}$

4 $r_n \sim 1.27 \frac{1.81^n}{n\sqrt{n}}$

Objective: Count π' -shapes compatible with RNA of length n .

1 π' -shapes = bracket words avoiding motifs $[[\dots]]$ and $\bullet\bullet$

2 $R \rightarrow \square R \mid S \quad S \rightarrow U[T]S \mid U \quad U \rightarrow \diamond \mid \varepsilon$

$T \rightarrow U[T]U[T]S \mid \diamond[T] \mid [T]\diamond \mid \diamond[T]\diamond \mid \bullet^\theta$

3 $\theta = 3, R(z) = \frac{1+2z^2+2z^3+z^4-z^5-z^6-\sqrt{1-4z^3-2z^4-2z^5+2z^6-7z^8-z^{10}+2z^{11}+z^{12}}}{2z^2(1-z^2)}$

4 $r_n \sim 1.27 \frac{1.81^n}{n\sqrt{n}}$

Objective: Count π' -shapes compatible with RNA of length n .

1 π' -shapes = bracket words avoiding motifs $[[\dots]]$ and $\bullet\bullet$

2 $R \rightarrow \square R \mid S \quad S \rightarrow U[T]S \mid U \quad U \rightarrow \diamond \mid \varepsilon$

$T \rightarrow U[T]U[T]S \mid \diamond[T] \mid [T]\diamond \mid \diamond[T]\diamond \mid \bullet^\theta$

3 $\theta = 3, R(z) = \frac{1+2z^2+2z^3+z^4-z^5-z^6-\sqrt{1-4z^3-2z^4-2z^5+2z^6-7z^8-z^{10}+2z^{11}+z^{12}}}{2z^2(1-z^2)}$

4 $r_n \sim 1.27 \frac{1.81^n}{n\sqrt{n}}$

Objective: Count π' -shapes compatible with RNA of length n .

1 π' -shapes = bracket words avoiding motifs $[[\dots]]$ and $\bullet\bullet$

2 $R \rightarrow \square R \mid S \quad S \rightarrow U[T]S \mid U \quad U \rightarrow \diamond \mid \varepsilon$

$T \rightarrow U[T]U[T]S \mid \diamond[T] \mid [T]\diamond \mid \diamond[T]\diamond \mid \bullet^\theta$

3 $\theta = 3, R(z) = \frac{1+2z^2+2z^3+z^4-z^5-z^6-\sqrt{1-4z^3-2z^4-2z^5+2z^6-7z^8-z^{10}+2z^{11}+z^{12}}}{2z^2(1-z^2)}$

4 $r_n \sim 1.27 \frac{1.81^n}{n\sqrt{n}}$

Model	Asymptotic number
Sec. str. on n – Combinatorial	$1.1 \cdot \frac{2.6^n}{n\sqrt{n}}$
Sec. str. on n – Empirical	$0.04 \cdot \frac{1.4^n}{n\sqrt{n}}$
π -shapes of size n	$1.38 \cdot \frac{1.73^n}{n\sqrt{n}}$
π -shapes compatible with sec. str. on n	$2.44 \cdot \frac{1.32^n}{n\sqrt{n}}$
π -shapes – Empirical	$0.21 \cdot \frac{1.1^n}{n\sqrt{n}}$
π' -shapes of size n	$0.99 \cdot \frac{2.41^n}{n\sqrt{n}}$
π' -shapes compatible with sec. str. on n	$1.28 \cdot \frac{1.81^n}{n\sqrt{n}}$

Message #4

The number of π -shapes (coarsest) is upper-bounded by 1.3^n .
 ⇒ Shape-based indexing is feasible through exhaustive search.

Part III

Realistic models for random RNA structures

Random generation can be used to:

- Assess the significance of observed phenomena
- Estimate the practical complexity of algorithms

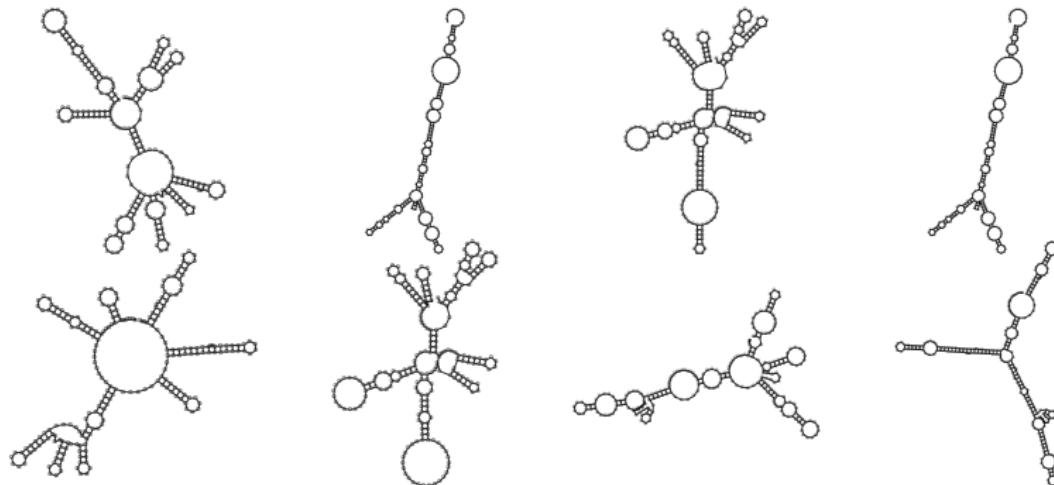
Worst-case complexity can be *misleading*:

- Statistical sampling: $\Theta(n^2) \rightarrow \mathcal{O}(n\sqrt{n})$
- Tree alignment: $\Theta(n^4) \rightarrow \mathcal{O}(n^2)$

Average case analysis: Hidden uniformity hypothesis

Yet all secondary structures are not equally likely observed as native structures!

RNA secondary structures (Mathews DB)



Grammar approach for the random generation

RNA secondary structures = $()$ -free Motzkin words [Wat78]

+ Uniform random generation (Boltzmann, recursive ...)

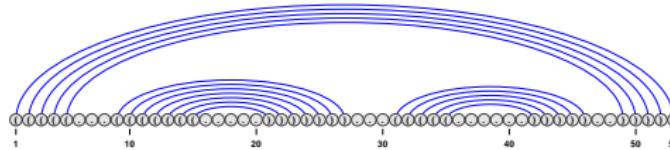
Grammar approach for the random generation

RNA secondary structures = $()$ -free Motzkin words [Wat78]

+ Uniform random generation (Boltzmann, recursive ...)

Example:

$(((((\bullet\bullet\bullet(((((\bullet\bullet\bullet\bullet)))))))\bullet\bullet\bullet((((\bullet\bullet\bullet\bullet))))\bullet\bullet))))$



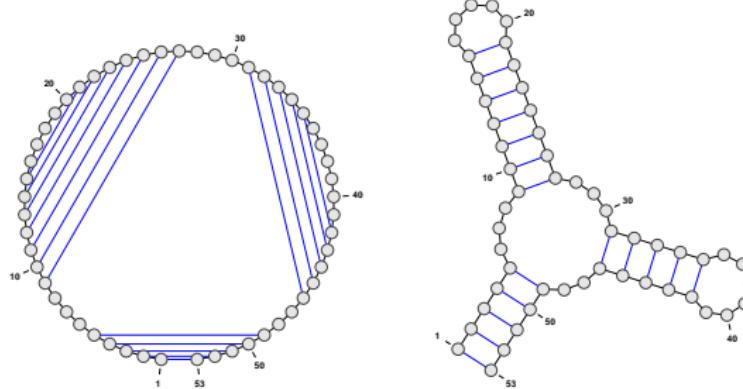
Grammar approach for the random generation

RNA secondary structures = $()$ -free Motzkin words [Wat78]

+ Uniform random generation (Boltzmann, recursive ...)

Example:

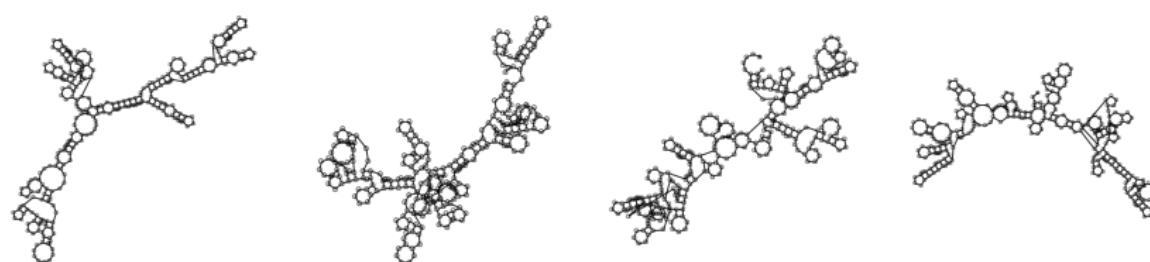
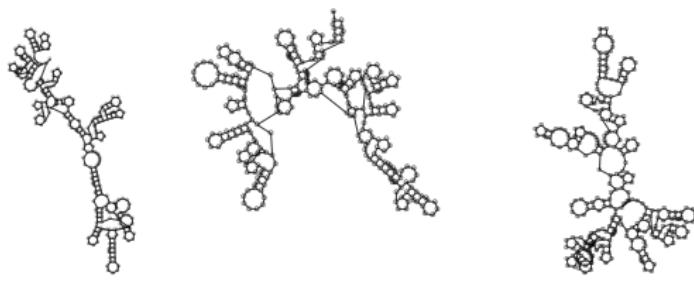
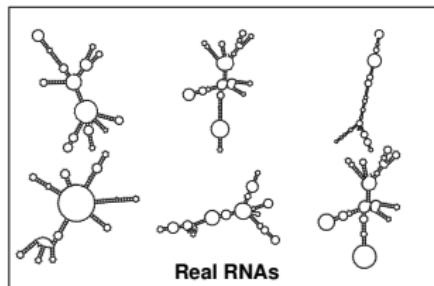
$(((((\bullet\bullet\bullet(((((\bullet\bullet\bullet\bullet)))))))\bullet\bullet\bullet((((\bullet\bullet\bullet\bullet))))\bullet\bullet))))$



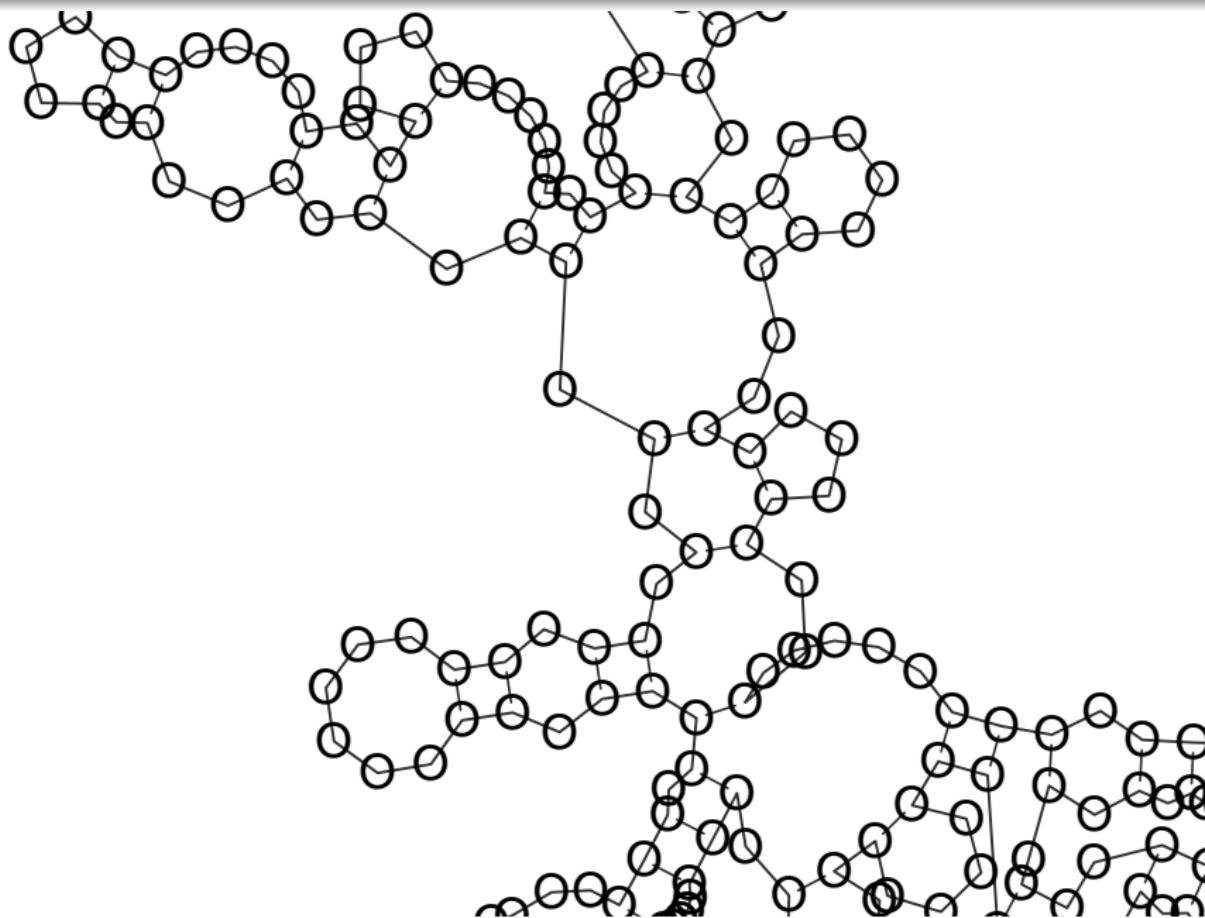
Grammar approach for the random generation

RNA secondary structures = $()$ -free Motzkin words [Wat78]

+ Uniform random generation (Boltzmann, recursive ...)



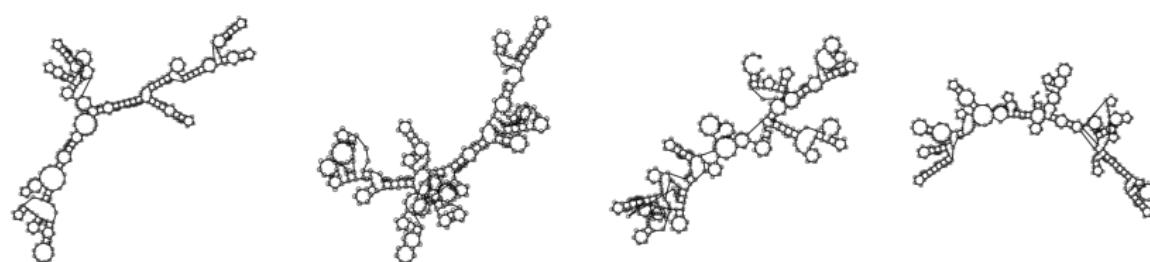
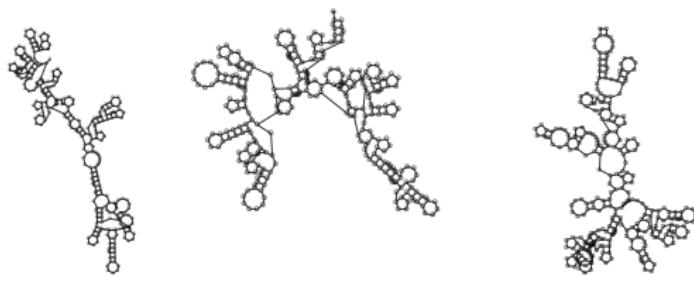
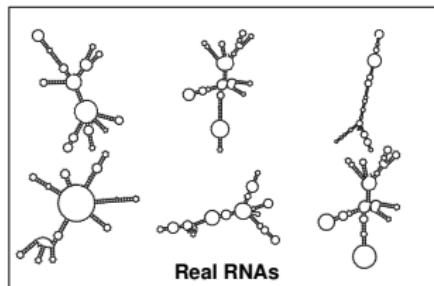
Grammar approach for the random generation



Grammar approach for the random generation

RNA secondary structures = $()$ -free Motzkin words [Wat78]

+ Uniform random generation (Boltzmann, recursive ...)



RNA secondary structures = (\cdot) -free Motzkin words [Wat78]

- + Uniform random generation (Boltzmann, recursive ...)
- + Sequence/structure annotation (h : Base-pair, H : Helix, ...)

Example:

$(((((\bullet\bullet\bullet((((((\bullet\bullet\bullet\bullet))))))))\bullet\bullet\bullet((((((\bullet\bullet\bullet\bullet))))))))\bullet\bullet))))$
 $\text{H}\text{h}\text{h}\text{h}\text{M}\text{m}\text{m}\text{H}\text{h}\text{h}\text{h}\text{h}\text{T}\text{t}\text{t}\text{t}\text{t}\text{h}\text{h}\text{h}\text{h}\text{h}\text{m}\text{m}\text{m}\text{H}\text{h}\text{h}\text{h}\text{T}\text{t}\text{t}\text{t}\text{t}\text{t}\text{h}\text{h}\text{h}\text{h}\text{m}\text{m}\text{m}\text{h}\text{h}\text{h}\text{h}$

Annotate native structures from Mathews database
⇒ Proportions of occurrences of each terminal

Grammar can be changed to generate already annotated structures
⇒ Evaluate expected proportions in the uniform model

Available techniques:

- Empirical
- Recursive/Dynamic programming
- Flajolet–Symbolic, ...

Grammar approach for the random generation

RNA secondary structures = (\cdot) -free Motzkin words [Wat78]

- + Uniform random generation (Boltzmann, recursive ...)
- + Sequence/structure annotation (h : Base-pair, H : Helix, ...)

Feature	B	b	I	i	M	m	T	t	H	h
\mathcal{M}_0 (%)	7.2	5.6	2.8	7.3	3.7	7.6	5.2	14.5	18.6	27.5
Exper.	1.5	2.3	1.9	11.2	1.1	9.0	2.6	16.6	4.8	48.9

Proportions of symbols associated with structural features:
Uniform model \mathcal{M}_0 vs Observed ($n = 300$)

As suspected, avg helix length $(\text{H} + \text{h})/(2 * \text{H})$ is way off!
(\mathcal{M}_0 : 1.23, Exp: 5.59)

Message #1

RNA secondary structures \Rightarrow Uniform model = Bad!

Null-hypothesis too easily refuted/Overrated statistical significance
Bonzaï structures might yield biased estimates.

Message #2

Need (at least) to account for constraints on **min/avg** sizes of features.

- Helices
- Sequences of unpaired bases
- Degree of branching loops

Constraints captured by grammars in two ways:

- **Min:** Tweak grammar
- **Avg:** Modify distribution \Rightarrow Weighted probability distribution

Definition (Context-free grammar)

Context-free grammar = 4-tuple $(\Sigma, \mathcal{N}, \mathcal{P}, S)$:

- Σ : Alphabet.
- \mathcal{N} : Non-terminal symbols.
- \mathcal{P} : Set of production rules $N \rightarrow X \in \mathcal{N} \times \{\Sigma \cup \mathcal{N}\}^*$.
- S : Axiom, or initial non-terminal.

Definition (**Weighted** context-free grammar [DRT00])

A **weighted** context-free grammar is a **5-tuple** $\mathcal{G}_\pi = (\Sigma, \mathcal{N}, \mathcal{P}, \mathcal{S}, \pi)$:

- $\Sigma, \mathcal{N}, \mathcal{P}, \mathcal{S}$: Same as previously.
- π : **Weight function** $\pi : \Sigma \rightarrow \mathbb{R}$.

Then \mathcal{G}_π implicitly defines a **weighted probability distribution** \mathcal{W} :

$$\forall \omega \in \mathcal{L}(\mathcal{G}_\pi), \quad \mathbb{P}(\omega) = \frac{\pi(\omega)}{\pi(\mathcal{L}(\mathcal{G}_\pi))} = \frac{\prod_{i=1}^{|\omega|} \pi(\omega_i)}{\pi(\mathcal{L}(\mathcal{G}_\pi))}.$$

Generation: Generating k words of size n takes $\mathcal{O}(n^2 + n \log(n).k)^*$ or $\mathcal{O}(n \log(n).k)^*$ using symbolic preprocessing (linear recurrences).

* Arithmetic operations $\Rightarrow \Theta(n)$ overhead OR tricky interval arithmetics [DZ99]

Definition (**Weighted** context-free grammar [DRT00])

A **weighted** context-free grammar is a **5-tuple** $\mathcal{G}_\pi = (\Sigma, \mathcal{N}, \mathcal{P}, \mathcal{S}, \pi)$:

- $\Sigma, \mathcal{N}, \mathcal{P}, \mathcal{S}$: Same as previously.
- π : **Weight function** $\pi : \Sigma \rightarrow \mathbb{R}$.

Then \mathcal{G}_π implicitly defines a **weighted probability distribution** \mathcal{W} :

$$\forall \omega \in \mathcal{L}(\mathcal{G}_\pi), \quad \mathbb{P}(\omega) = \frac{\pi(\omega)}{\pi(\mathcal{L}(\mathcal{G}_\pi))} = \frac{\prod_{i=1}^{|\omega|} \pi(\omega_i)}{\pi(\mathcal{L}(\mathcal{G}_\pi))}.$$

Generation: Generating k words of size n takes $\mathcal{O}(n^2 + n \log(n).k)^*$ or $\mathcal{O}(n \log(n).k)^*$ using symbolic preprocessing (linear recurrences).

* Arithmetic operations $\Rightarrow \Theta(n)$ overhead OR tricky interval arithmetics [DZ99]

1 – Uniform models are bad!

2 – Let's use weighted distribution!

Which one?

Weights should be computed such that:

- Sampled structure follow Boltzmann distribution
Not very good for modeling native structure
- Expectation is maximized
Not very easy ($\text{WCFG} \neq \text{SCFG}$)
Overfitting vs Inconsistency
- Aiming at observed average properties
Average random objects should mimic features of the training set

How to compute weights targeting a given observation?

1 – Uniform models are bad!

2 – Let's use weighted distribution!

Which one?

Weights should be computed such that:

- Sampled structure follow Boltzmann distribution
Not very good for modeling native structure
- Expectation is maximized
Not very easy ($\text{WCFG} \neq \text{SCFG}$)
Overfitting vs Inconsistency
- Aiming at observed average properties
Average random objects should mimic features of the training set

How to compute weights targeting a given observation?

1 – Uniform models are bad!

2 – Let's use weighted distribution!

Which one?

Weights should be computed such that:

- Sampled structure follow Boltzmann distribution
Not very good for modeling native structure
- Expectation is maximized
Not very easy ($\text{WCFG} \neq \text{SCFG}$)
Overfitting vs Inconsistency
- Aiming at observed average properties
Average random objects should mimic features of the training set

How to compute weights targeting a given observation?

$f_{n,i}^\pi$: Proportion of Z_i in words of size n drawn from a π -distribution

$f_{*,i}^\pi$: Asymptotic proportion of Z_i in π ($f_{*,i}^\pi = \lim_{n \rightarrow \infty} f_{n,i}^\pi$)

Problem

Input: Grammar \mathcal{G} , target frequencies (μ_i) , length n

Goal: Find π_0 such that $\forall i \in [1, |\Sigma|]$, $f_{n,i}^{\pi_0} = \mu_i$

- Analytic (asymptotic) Rational
- Analytic (asymptotic) Algebraic
- Automatic heuristic (Fixed size)

Computation of $f_{n,i}^\pi$ easier than that of weight vector π_0

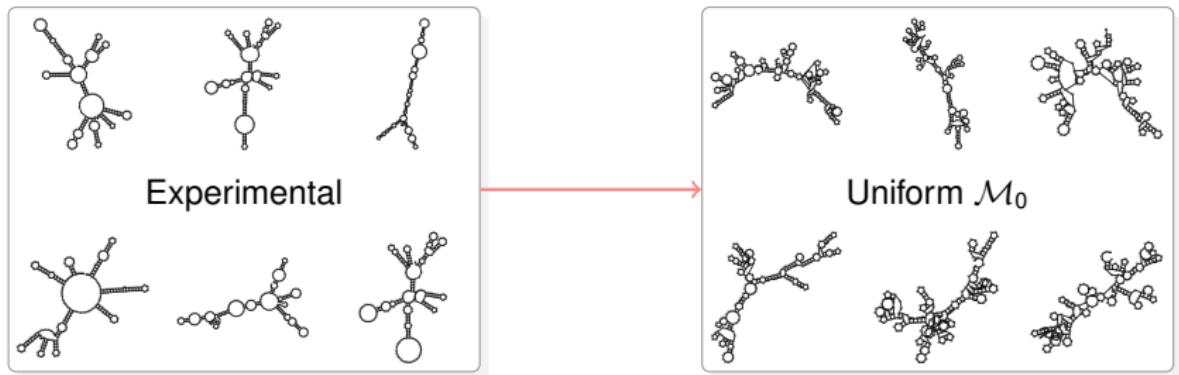
+ Partial monotonicity: $(\pi(Z_i) \nearrow) \rightarrow (f_{n,i}^\pi \nearrow)$

⇒ Rephrase as an continuous optimization problem

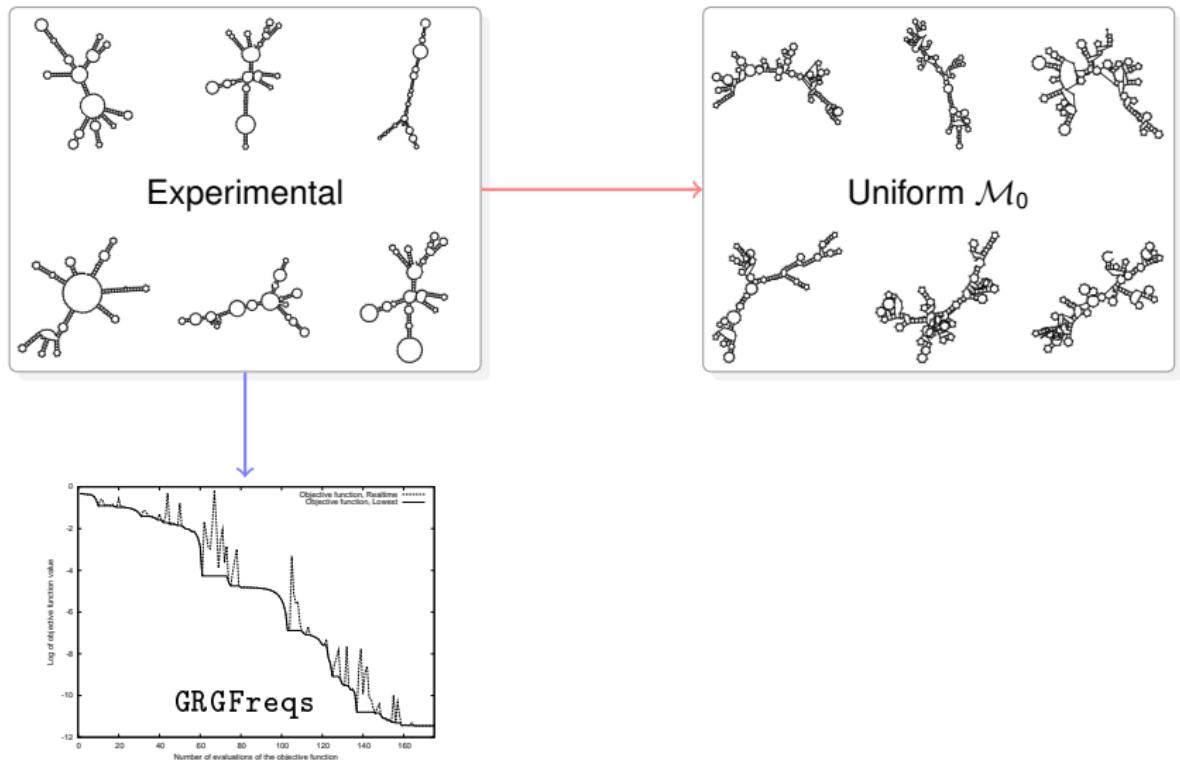
$$\text{Objective function: } g(\pi) = \sqrt{\sum_i (f_{n,i}^\pi - \mu_i)^2}$$

We use the CONDOR optimizer software package [Ber05].

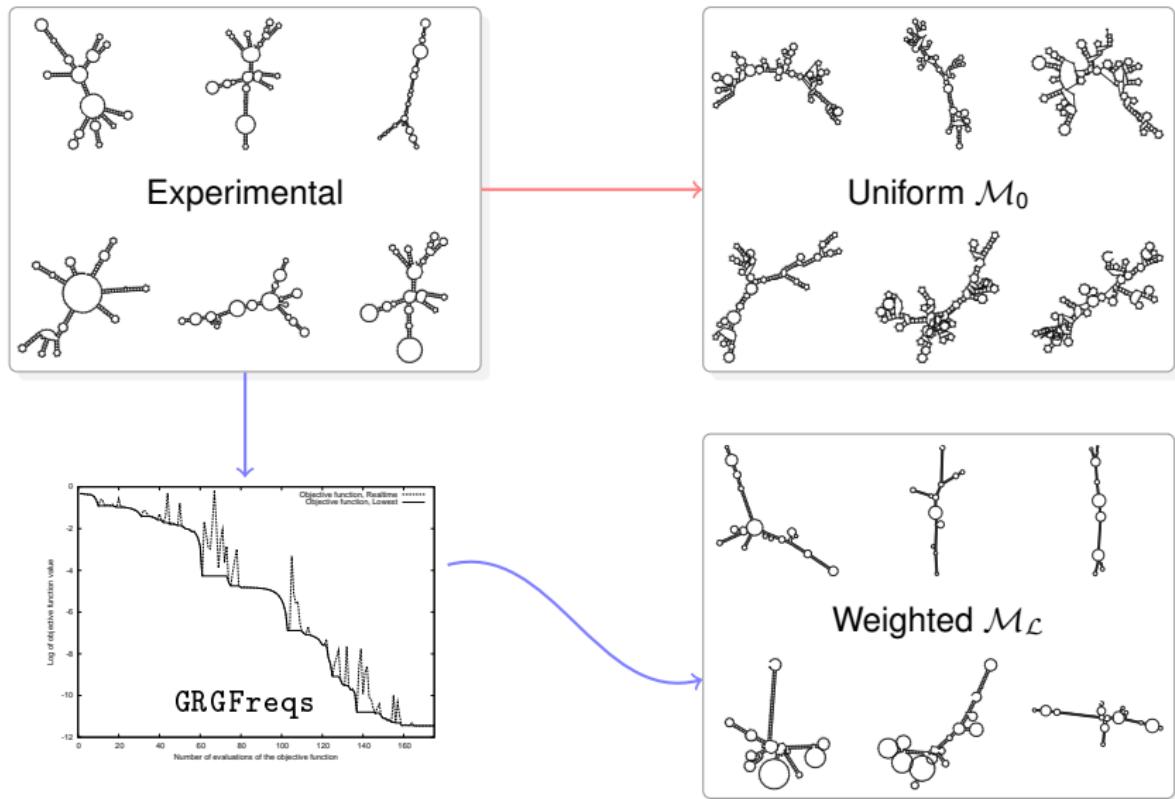
Application



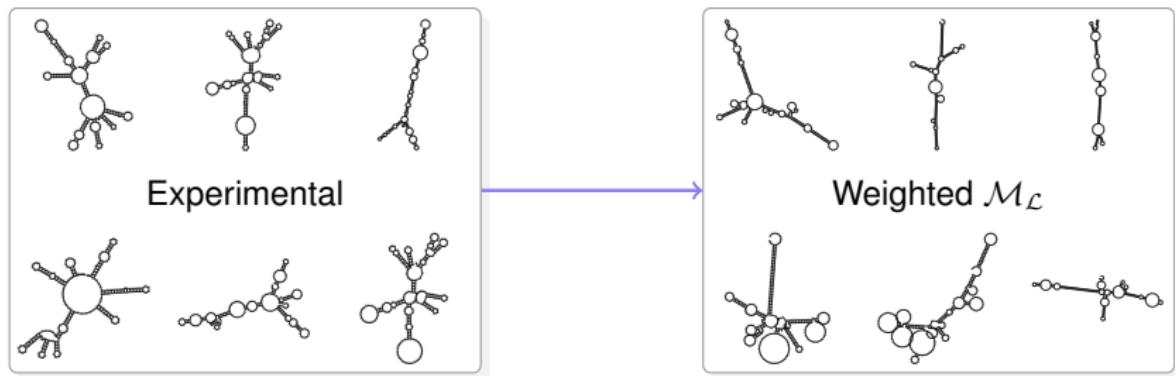
Application



Application

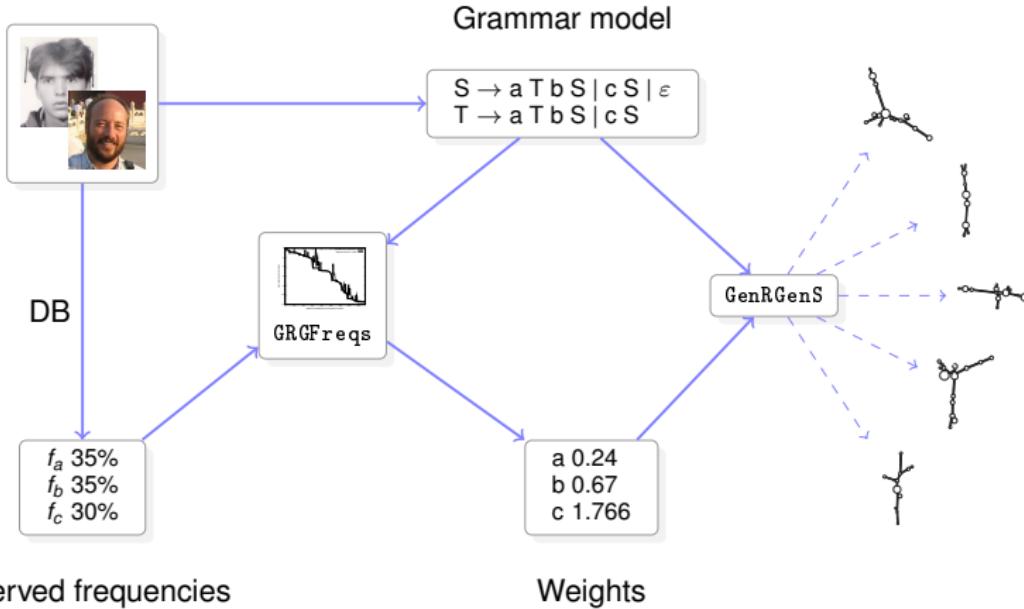


Application



Feature	B	b	I	i	M	m	T	t	H	h
Weight	1.15	1.15	1.15	1.15	2.75	1	1.15	1.15	.00308	1.25
\mathcal{M}_0 (%)	7.2	5.6	2.8	7.3	3.7	7.6	5.2	14.5	18.6	27.5
\mathcal{M}_L (%)	0.6	2.9	1.5	16.1	1.1	9.0	1.8	13.3	4.8	48.9
Exp.	1.5	2.3	1.9	11.2	1.1	9.0	2.6	16.6	4.8	48.9

People

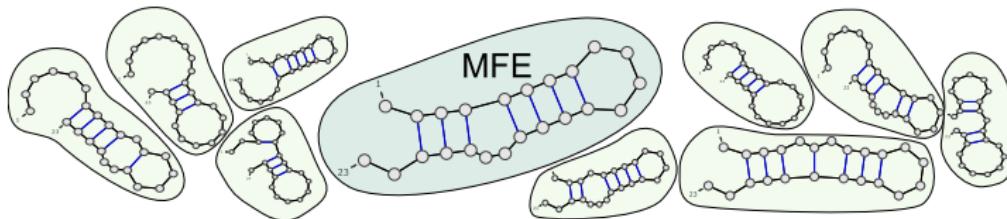


Now we can generate random structures fulfilling any desired properties.

Part IV

Improved statistical sampling

RNA *breathes* \Rightarrow There is not necessarily one functional conformation!



Boltzmann-ensemble paradigm

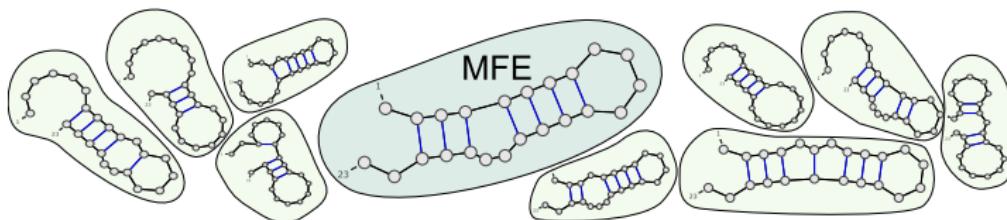
All possible secondary structures for an RNA coexist in a probability distribution, the **Boltzmann distribution**.

Consequently: MFE structure probability can be significantly remote (*unstructured RNAs?*).

Conversely, base-pairs having probability $> .99$ are valid 95% of the time.

- ⇒ Similar structures can **bundle up**.
- ⇒ Functional structures must be sought in **sub-optimal structures**.

RNA *breathes* \Rightarrow There is not necessarily one functional conformation!



Boltzmann-ensemble paradigm

All possible secondary structures for an RNA coexist in a probability distribution, the Boltzmann distribution.

Experiment: [DCL05]

- Sample structures from Boltzmann probability
 - Cluster based on structural distance
 - Build consensus structure based on most-probable cluster
- \Rightarrow Relative improvement for specificity (+17.6%) and sensitivity (+21.74%, except Group II Introns)

Boltzmann distribution: Definitions

A Boltzmann distribution assigns a weight to each structure S for an RNA ω , called the Boltzmann factor $e^{\frac{-E_{S,\omega}}{RT}}$ where:

- $E_{S,\omega}$ is the free-energy of S (kCal.mol^{-1})
- T is the temperature (K)
- R is the Boltzmann constant ($1.986 \cdot 10^{-3} \text{ kCal.K}^{-1} \cdot \text{mol}^{-1}$)

Distribution is renormalized on \mathcal{S}_ω by the partition function

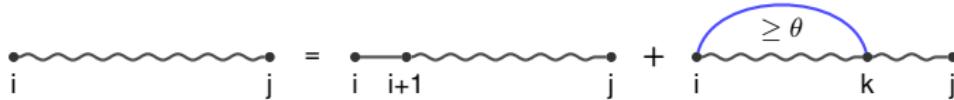
$$\mathcal{Z}_\omega = \sum_{S \in \mathcal{S}_\omega} e^{\frac{-E_{S,\omega}}{RT}}.$$

where \mathcal{S}_ω is the set of conformations compatible with ω .

The Boltzmann probability of a structure S is therefore

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

Partition function



Recurrence on the **minimal free-energy** of a folding:

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

$$N_{i,j} = \min \left\{ \begin{array}{ll} N_{i+1,j} & (i \text{ unpaired}) \\ \min_{k=i+\theta+1}^j E_{i,j} + N_{i+1,k-1} + N_{k+1,j} & (i \text{ comp. with } k) \end{array} \right.$$

+ Unambiguity

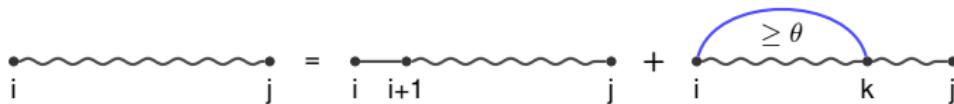
⇒ Recurrence on **partition function**:

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

$$\mathcal{Z}_{i,j} = \sum \left\{ \begin{array}{ll} \mathcal{Z}_{i+1,j} & (i \text{ unpaired}) \\ \sum_{k=i+\theta+1}^j e^{\frac{1}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} & (i \text{ comp. with } k) \end{array} \right.$$

+ Some slightly non-trivial computations...

= **Base-pair probabilities** in the Boltzmann ensemble [McC90].



Summands of $\mathcal{Z} \Leftrightarrow$ Cumulated Boltzmann weight of all accessible structures.

\Rightarrow Stochastic Backtrack recursively generates compatible structures with respect to the Boltzmann distribution [DL03].

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

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

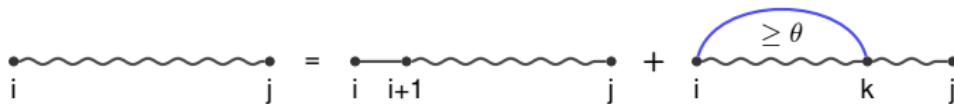
A

B_k

$A | B_{i+\theta+1} | B_{i+\theta+2} | B_{i+\theta+3} | \dots | B_{j-2} | B_{j-1} | B_j$

\Rightarrow Worst-case complexity in $\mathcal{O}(n^2 k)$ for k samples.

Could we do better than that to sample more at the same computational cost?



Summands of $\mathcal{Z} \Leftrightarrow$ Cumulated Boltzmann weight of all accessible structures.

\Rightarrow Stochastic Backtrack recursively generates compatible structures with respect to the Boltzmann distribution [DL03].

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

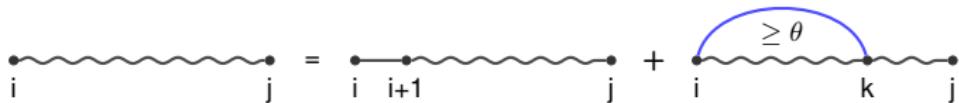
$$\mathcal{Z}_{i,j} = \sum_{k=i+\theta+1}^j e^{\frac{1}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j}$$

Diagram illustrating the recursive formula for $\mathcal{Z}_{i,j}$. The term $\mathcal{Z}_{i,j}$ is shown as a sum of terms $\mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j}$ for k from $i+\theta+1$ to j . Red annotations show a red bracket under the first term $\mathcal{Z}_{i+1,k-1}$ with a red arrow pointing to it, and another red bracket under the entire sum with a red arrow pointing to it. Red boxes labeled '???' are placed around the first term and the entire sum. To the right, there are two circles: one blue circle labeled 'A' and one pink circle labeled 'B_k'.

$$A | B_{i+\theta+1} | B_{i+\theta+2} | B_{i+\theta+3} | \dots | B_{j-2} | B_{j-1} | B_j$$

\Rightarrow Worst-case complexity in $\mathcal{O}(n^2 k)$ for k samples.

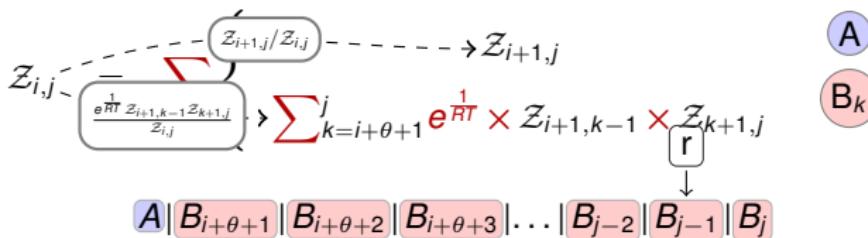
Could we do better than that to sample more at the same computational cost?



Summands of $\mathcal{Z} \Leftrightarrow$ Cumulated Boltzmann weight of all accessible structures.

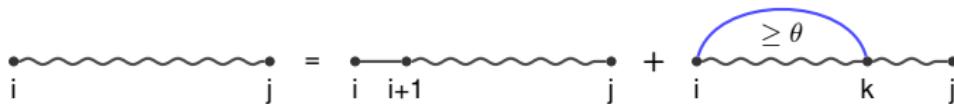
\Rightarrow Stochastic Backtrack recursively generates compatible structures with respect to the Boltzmann distribution [DL03].

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



\Rightarrow Worst-case complexity in $\mathcal{O}(n^2 k)$ for k samples.

Could we do better than that to sample more at the same computational cost?



Summands of $\mathcal{Z} \Leftrightarrow$ Cumulated Boltzmann weight of all accessible structures.

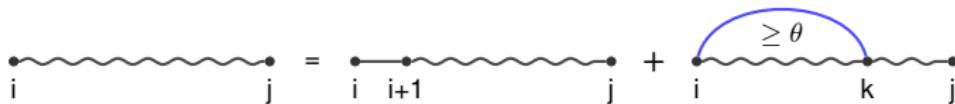
\Rightarrow Stochastic Backtrack recursively generates compatible structures with respect to the Boltzmann distribution [DL03].

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

$$\mathcal{Z}_{i,j} = \sum \left\{ \begin{array}{c} \mathcal{Z}_{i+1,j} \\ \sum_{k=i+\theta+1}^j e^{\frac{1}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \end{array} \right. \begin{array}{l} \text{A} \\ \text{B}_k \\ r \\ \downarrow \\ A | B_{i+\theta+1} | B_{i+\theta+2} | B_{i+\theta+3} | \dots | B_{j-2} | B_{j-1} | B_j \end{array}$$

\Rightarrow Worst-case complexity in $\mathcal{O}(n^2 k)$ for k samples.

Could we do better than that to sample more at the same computational cost?



Summands of $\mathcal{Z} \Leftrightarrow$ Cumulated Boltzmann weight of all accessible structures.

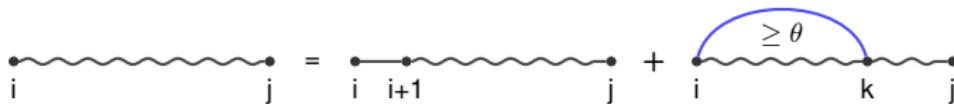
\Rightarrow Stochastic Backtrack recursively generates compatible structures with respect to the Boltzmann distribution [DL03].

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

$$\mathcal{Z}_{i,j} = \sum \left\{ \begin{array}{c} \mathcal{Z}_{i+1,j} \\ \sum_{k=i+\theta+1}^j e^{\frac{1}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \end{array} \right. \begin{array}{l} \text{A} \\ \text{B}_k \\ r \\ \downarrow \\ A | B_{i+\theta+1} | B_{i+\theta+2} | B_{i+\theta+3} | \dots | B_{j-2} | B_{j-1} | B_j \end{array}$$

\Rightarrow Worst-case complexity in $\mathcal{O}(n^2 k)$ for k samples.

Could we do better than that to sample more at the same computational cost?



Summands of $\mathcal{Z} \Leftrightarrow$ Cumulated Boltzmann weight of all accessible structures.

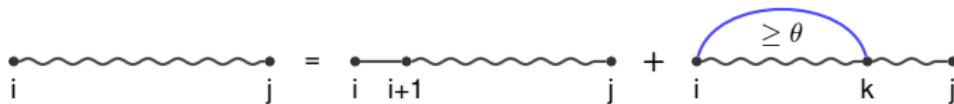
\Rightarrow Stochastic Backtrack recursively generates compatible structures with respect to the Boltzmann distribution [DL03].

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

$$\mathcal{Z}_{i,j} = \sum \left\{ \begin{array}{c} \mathcal{Z}_{i+1,j} \\ \sum_{k=i+\theta+1}^j e^{\frac{1}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \end{array} \right. \begin{array}{l} \text{A} \\ \text{B}_k \\ r \\ \downarrow \\ A | B_{i+\theta+1} | B_{i+\theta+2} | B_{i+\theta+3} | \dots | B_{j-2} | B_{j-1} | B_j \end{array}$$

\Rightarrow Worst-case complexity in $\mathcal{O}(n^2 k)$ for k samples.

Could we do better than that to sample more at the same computational cost?



Summands of $\mathcal{Z} \Leftrightarrow$ Cumulated Boltzmann weight of all accessible structures.

\Rightarrow Stochastic Backtrack recursively generates compatible structures with respect to the Boltzmann distribution [DL03].

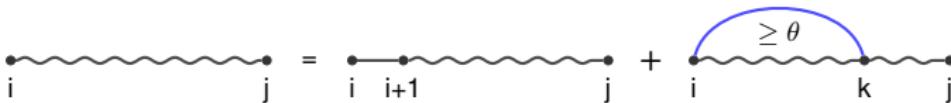
$$\mathcal{Z}_{i,t} = 1, \quad \forall t \in [i, i + \theta]$$

$$\mathcal{Z}_{i,j} = \sum \left\{ \begin{array}{c} \mathcal{Z}_{i+1,j} \\ \sum_{k=i+\theta+1}^j e^{\frac{1}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} \end{array} \right. \begin{array}{l} \text{A} \\ \text{B}_k \\ r \end{array}$$

Below the equation, a sequence of RNA segments is shown: $A | B_{i+\theta+1} | B_{i+\theta+2} | B_{i+\theta+3} | \dots | B_{j-2} | B_{j-1} | B_j$. The segments $B_{i+\theta+1}, B_{i+\theta+2}, B_{i+\theta+3}, \dots, B_{j-2}, B_{j-1}$ are highlighted in pink. A dashed arrow points from the term $\mathcal{Z}_{k+1,j}$ in the equation to the segment B_j in the sequence. A small circle labeled r is placed above the segment B_j .

\Rightarrow Worst-case complexity in $\mathcal{O}(n^2 k)$ for k samples.

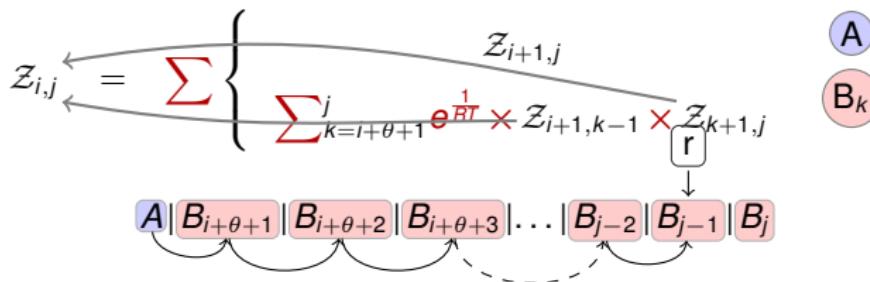
Could we do better than that to sample more at the same computational cost?



Summands of $\mathcal{Z} \Leftrightarrow$ Cumulated Boltzmann weight of all accessible structures.

\Rightarrow Stochastic Backtrack recursively generates compatible structures with respect to the Boltzmann distribution [DL03].

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



\Rightarrow Worst-case complexity in $\mathcal{O}(n^2k)$ for k samples.

Could we do better than that to sample more at the same computational cost?

Average-case analysis

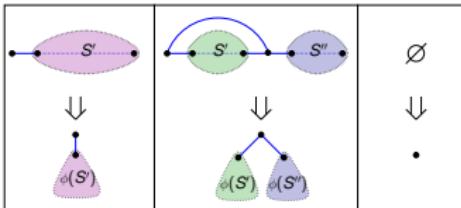
Homopolymer model: Any base-pair allowed.

Boltzmann distribution: Based on a Nussinov/Jacobson energy model.

Theorem

Let n be the RNA length et k the #samples.

The average-case complexity of statistical sampling is in $\Theta(n^3 + kn\sqrt{n})$.



Remark 1: Cost $c(S)$ of generating a structure S is dominated by the cumulated size of the leftward children $t \in \mathcal{T}_g(S)$ in the associated unary/binary tree:

$$c(S) = n - \#bp(S)/2 + \sum_{t \in \mathcal{T}_g(S)} |t| - \theta$$

Average-case analysis

Homopolymer model: Any base-pair allowed.

Boltzmann distribution: Based on a Nussinov/Jacobson energy model.

Theorem

Let n be the RNA length et k the #samples.

The average-case complexity of statistical sampling is in $\Theta(n^3 + kn\sqrt{n})$.

Let us focus on the expectation of the *generation cost* random variable X . We have

$$\mathbb{E}(X \mid n) = \sum_{|S|=n} c(S) \cdot \frac{e^{\frac{bp(S)}{RT}}}{Z_n} = \frac{\sum_{|S|=n} c(S) \cdot e^{\frac{bp(S)}{RT}}}{\sum_{|S|=n} e^{\frac{bp(S)}{RT}}}.$$

Let us consider separately the two generating functions $C(z)$ and $P_f(z)$ such that

$$C(z) = \sum_{S \in \mathcal{S}} e^{\frac{bp(S)}{RT}} c(S) z^{|S|} \quad \text{et} \quad P_f(z) = \sum_{S \in \mathcal{S}} e^{\frac{bp(S)}{RT}} z^{|S|}.$$

Let us remark that $\mathbb{E}(X \mid n) = [z^n]C(z)/[z^n]P_f(z)$.

Average-case analysis

Homopolymer model: Any base-pair allowed.

Boltzmann distribution: Based on a Nussinov/Jacobson energy model.

Theorem

Let n be the RNA length et k the #samples.

The average-case complexity of statistical sampling is in $\Theta(n^3 + kn\sqrt{n})$.

$C(z)$ and $P_f(z)$ are the positive solutions of the system

$$\begin{aligned} C(z) &= z(P_f(z) + C(z)) + z^2 e^{\frac{1}{RT}} (1 - \theta) P_f^{\geq \theta}(z) P_f(z) \\ &+ z^3 e^{\frac{1}{RT}} \frac{\partial P_f^{\geq \theta}(z)}{\partial z} P_f(z) + z^2 e^{\frac{1}{RT}} C^{\geq \theta}(z) P_f(z) \\ &+ z^2 e^{\frac{1}{RT}} P_f^{\geq \theta}(z) C(z) \\ P_f(z) &= z^2 e^{\frac{1}{RT}} P_f^{\geq \theta}(z) P_f(z) + z P_f(z) + 1 \end{aligned}$$

where $S^{\geq \theta}(z)$ is the truncated gen. fun. $S(z)$, deprived of its terms of degrees $< \theta$.

Average-case analysis

Homopolymer model: Any base-pair allowed.

Boltzmann distribution: Based on a Nussinov/Jacobson energy model.

Theorem

Let n be the RNA length et k the #samples.

The average-case complexity of statistical sampling is in $\Theta(n^3 + kn\sqrt{n})$.

⇒ Big generating functions (Thanks to Maple !)...

+ Singularity analysis:

$$[z^n]P_f(z) \sim \frac{\kappa}{\rho^n n \sqrt{n}} (1 + \mathcal{O}(1/n)) \quad [z^n]C(z) \sim \frac{\kappa'}{\rho^n} (1 + \mathcal{O}(1/\sqrt{n}))$$

We obtain an asymptotic equivalent for $\mathbb{E}(X | n)$:

$$\mathbb{E}(X | n) = \frac{[z^n]C(z)}{[z^n]P_f(z)} \sim \frac{\kappa'}{\kappa} n \sqrt{n} (1 + \mathcal{O}(1/\sqrt{n})).$$

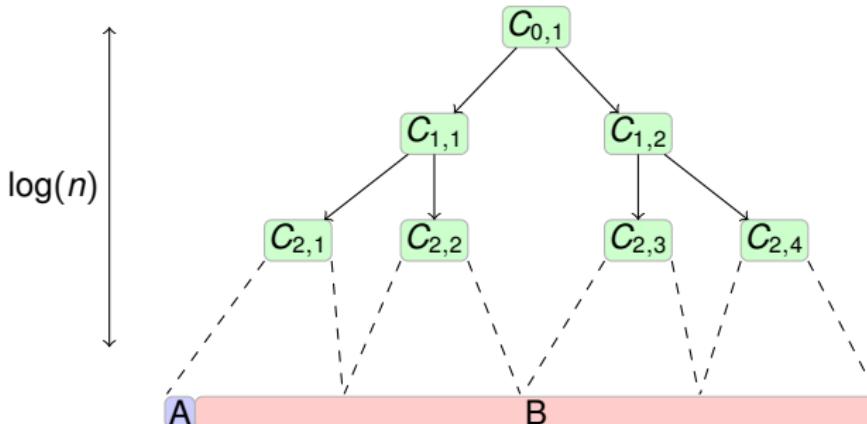
□

Remark: Bivariate version of the 4 steps program...

Boustrophédon search

Rationale: Most of the time is spent on testing unsuitable decompositions!

- Hierarchical organization of contributions



⇒ Generation in $\mathcal{O}(n^3 + kn \log(n))$, memory in $\Theta(n^3)$.

- Boustrophédon search



Unbalanced decompositions

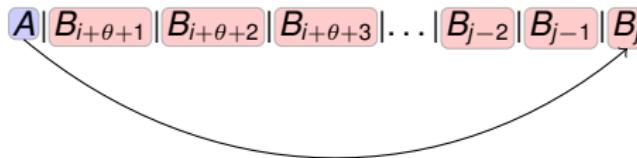
⇒ Low immediate cost, recursion on intervals of lengths $\mathcal{O}(c)/\mathcal{O}(j-i-c)$.

Rationale: Most of the time is spent on testing unsuitable decompositions!

- Hierarchical organization of contributions

⇒ Generation in $\mathcal{O}(n^3 + kn \log(n))$, memory in $\Theta(n^3)$.

- Boustrophédon search



Unbalanced decompositions

⇒ Low *immediate* cost, recursion on intervals of lengths $\mathcal{O}(c)/\mathcal{O}(j-i-c)$.

Décomposition égale

⇒ *Immediate* cost is high ($\mathcal{O}(j-i-c)$), divide and conquer.

Worst-case complexity solution of

$$f(n) = \max_{k \in [1, n-1]} (f(k) + f(n-k) + 2 \min(k, n-k)) \in \mathcal{O}(n \log(n)) \text{ [FZV94].}$$

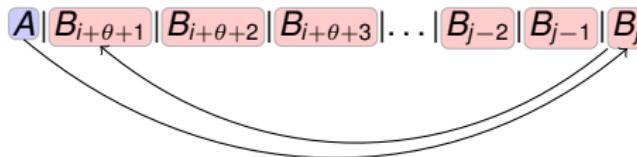
⇒ Generation in $\mathcal{O}(n^3 + kn \log(n))$ without any precomputation.

Rationale: Most of the time is spent on testing unsuitable decompositions!

- Hierarchical organization of contributions

⇒ Generation in $\mathcal{O}(n^3 + kn \log(n))$, memory in $\Theta(n^3)$.

- Boustrophédon search



Unbalanced decompositions

⇒ Low *immediate* cost, recursion on intervals of lengths $\mathcal{O}(c)/\mathcal{O}(j-i-c)$.

Décomposition égale

⇒ *Immediate* cost is high ($\mathcal{O}(j-i-c)$), divide and conquer.

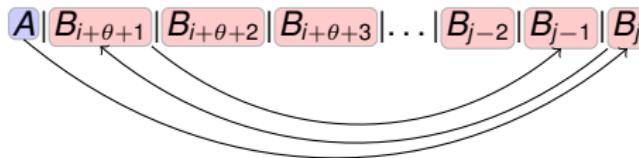
Worst-case complexity solution of

$$f(n) = \max_{k \in [1, n-1]} (f(k) + f(n-k) + 2 \min(k, n-k)) \in \mathcal{O}(n \log(n)) \text{ [FZV94].}$$

⇒ Generation in $\mathcal{O}(n^3 + kn \log(n))$ without any precomputation.

Rationale: Most of the time is spent on testing unsuitable decompositions!

- Hierarchical organization of contributions
⇒ Generation in $\mathcal{O}(n^3 + kn \log(n))$, memory in $\Theta(n^3)$.
- Boustrophédon search



Unbalanced decompositions

⇒ Low *immediate* cost, recursion on intervals of lengths $\mathcal{O}(c)/\mathcal{O}(j-i-c)$.

Décomposition égale

⇒ *Immediate* cost is high ($\mathcal{O}(j-i-c)$), divide and conquer.

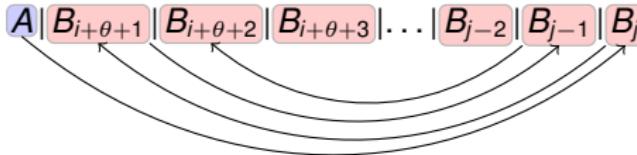
Worst-case complexity solution of

$$f(n) = \max_{k \in [1, n-1]} (f(k) + f(n-k) + 2 \min(k, n-k)) \in \mathcal{O}(n \log(n)) \text{ [FZV94].}$$

⇒ Generation in $\mathcal{O}(n^3 + kn \log(n))$ without any precomputation.

Rationale: Most of the time is spent on testing unsuitable decompositions!

- Hierarchical organization of contributions
⇒ Generation in $\mathcal{O}(n^3 + kn \log(n))$, memory in $\Theta(n^3)$.
- Boustrophédon search



Unbalanced decompositions

⇒ Low *immediate* cost, recursion on intervals of lengths $\mathcal{O}(c)/\mathcal{O}(j-i-c)$.

Décomposition égale

⇒ *Immediate* cost is high ($\mathcal{O}(j-i-c)$), divide and conquer.

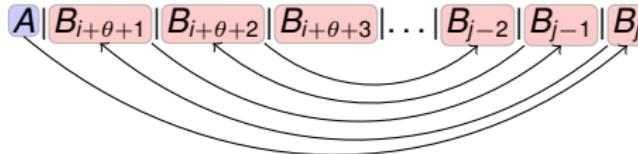
Worst-case complexity solution of

$$f(n) = \max_{k \in [1, n-1]} (f(k) + f(n-k) + 2 \min(k, n-k)) \in \mathcal{O}(n \log(n)) \text{ [FZV94].}$$

⇒ Generation in $\mathcal{O}(n^3 + kn \log(n))$ without any precomputation.

Rationale: Most of the time is spent on testing unsuitable decompositions!

- Hierarchical organization of contributions
⇒ Generation in $\mathcal{O}(n^3 + kn \log(n))$, memory in $\Theta(n^3)$.
- Boustrophédon search



Unbalanced decompositions

⇒ Low *immediate* cost, recursion on intervals of lengths $\mathcal{O}(c)/\mathcal{O}(j-i-c)$.

Décomposition égale

⇒ *Immediate* cost is high ($\mathcal{O}(j-i-c)$), divide and conquer.

Worst-case complexity solution of

$$f(n) = \max_{k \in [1, n-1]} (f(k) + f(n-k) + 2 \min(k, n-k)) \in \mathcal{O}(n \log(n)) \text{ [FZV94].}$$

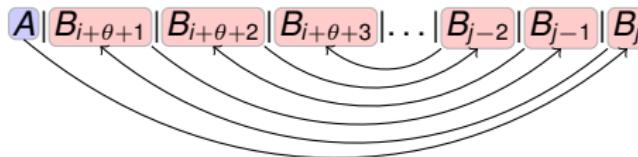
⇒ Generation in $\mathcal{O}(n^3 + kn \log(n))$ without any precomputation.

Rationale: Most of the time is spent on testing unsuitable decompositions!

- Hierarchical organization of contributions

⇒ Generation in $\mathcal{O}(n^3 + kn \log(n))$, memory in $\Theta(n^3)$.

- Boustrophédon search



Unbalanced decompositions

⇒ Low *immediate cost*, recursion on intervals of lengths $\mathcal{O}(c)/\mathcal{O}(j-i-c)$.

Décomposition égale

⇒ *Immediate cost is high* ($\mathcal{O}(j-i-c)$), divide and conquer.

Worst-case complexity solution of

$$f(n) = \max_{k \in [1, n-1]} (f(k) + f(n-k) + 2 \min(k, n-k)) \in \mathcal{O}(n \log(n)) \text{ [FZV94].}$$

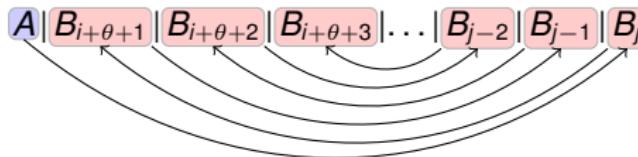
⇒ Generation in $\mathcal{O}(n^3 + kn \log(n))$ without any precomputation.

Rationale: Most of the time is spent on testing unsuitable decompositions!

- Hierarchical organization of contributions

⇒ Generation in $\mathcal{O}(n^3 + kn \log(n))$, memory in $\Theta(n^3)$.

- Boustrophédon search



Unbalanced decompositions

⇒ Low *immediate* cost, recursion on intervals of lengths $\mathcal{O}(c)/\mathcal{O}(j-i-c)$.

Décomposition égale

⇒ *Immediate* cost is high ($\mathcal{O}(j-i-c)$), divide and conquer.

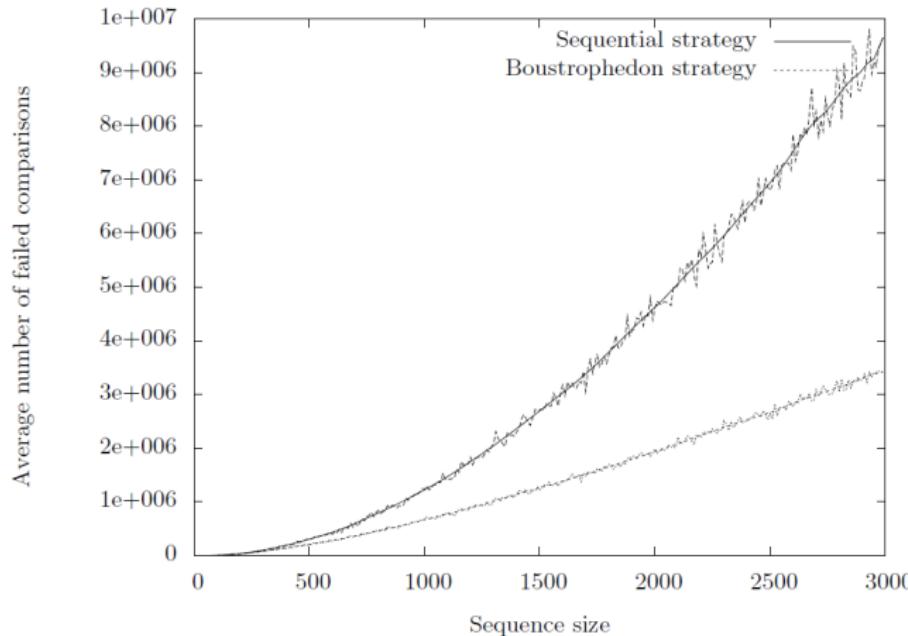
Worst-case complexity solution of

$$f(n) = \max_{k \in [1, n-1]} (f(k) + f(n-k) + 2 \min(k, n-k)) \in \mathcal{O}(n \log(n)) \text{ [FZV94].}$$

⇒ Generation in $\mathcal{O}(n^3 + kn \log(n))$ without any precomputation.

- Implementation in **Unafold** (successor of **MFold**) using **full Turner model**.
- Random sequences: Markov order 1 model based on a 5s rRNA.

Implementation in **Unafold** (successor of **MFold**).



Part V

Conclusion

Using enumerative (models) and analytic (asymptotics) techniques, one can:

- Accurately enumerate RNA Shapes [JCB08].



- Improved random models for structures [Bioinf.06, TCS10].
- Analyze and optimize statistical sampling algorithms [JMB08].



- Analyze Boltzmann samples sets of structures [GASCOM10].



- Traverse the full mutational landscape [Recomb11].



- Explore enriched conformational spaces (PKs, MC-Fold) [...].



- Address constrained RNA design problems [...].



Using enumerative (models) and analytic (asymptotics) techniques, one can:

- Accurately enumerate RNA Shapes [JCB08].



- Improved random models for structures [Bioinf.06, TCS10].



- Analyze and optimize statistical sampling algorithms [JMB08].



- Analyze Boltzmann samples sets of structures [GASCOM10].



- Traverse the full mutational landscape [Recomb11].



- Explore enriched conformational spaces (PKs, MC-Fold) [...].



- Address constrained RNA design problems [...].

Using enumerative (models) and analytic (asymptotics) techniques, one can:

- Accurately enumerate RNA Shapes [JCB08].
- Improved random models for structures [Bioinf.06, TCS10].
- Analyze and optimize statistical sampling algorithms [JMB08].
- Analyze Boltzmann samples sets of structures [GASCOM10].
- Traverse the full mutational landscape [Recomb11].
- Explore enriched conformational spaces (PKs, MC-Fold) [...].
- Address constrained RNA design problems [...].



Using enumerative (models) and analytic (asymptotics) techniques, one can:

- Accurately enumerate RNA Shapes [JCB08].
- Improved random models for structures [Bioinf.06, TCS10].
- Analyze and optimize statistical sampling algorithms [JMB08].
- Analyze Boltzmann samples sets of structures [GASCOM10].
- Traverse the full mutational landscape [Recomb11].
- Explore enriched conformational spaces (PKs, MC-Fold) [...].
- Address constrained RNA design problems [...].



Using enumerative (models) and analytic (asymptotics) techniques, one can:

- Accurately enumerate RNA Shapes [JCB08].
- Improved random models for structures [Bioinf.06, TCS10].
- Analyze and optimize statistical sampling algorithms [JMB08].
- Analyze Boltzmann samples sets of structures [GASCOM10].
- Traverse the full mutational landscape [Recomb11].
- Explore enriched conformational spaces (PKs, MC-Fold) [...].
- Address constrained RNA design problems [...].



Using enumerative (models) and analytic (asymptotics) techniques, one can:

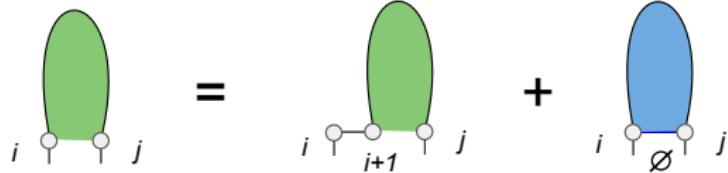
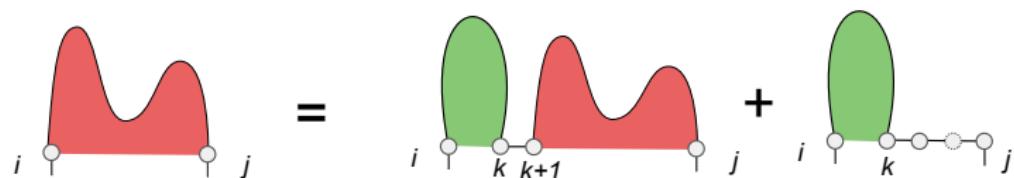
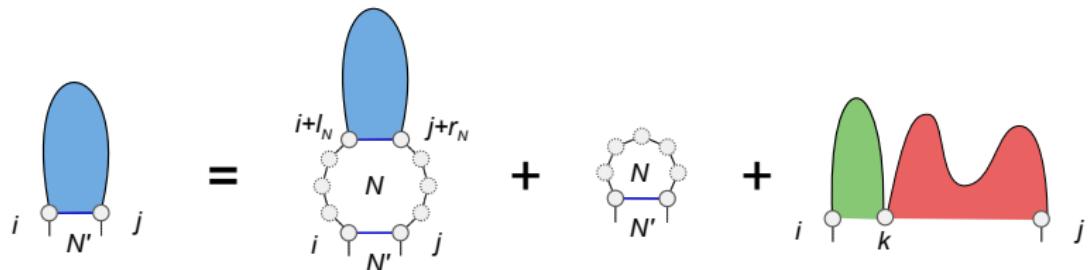
- Accurately enumerate RNA Shapes [JCB08].
- Improved random models for structures [Bioinf.06, TCS10].
- Analyze and optimize statistical sampling algorithms [JMB08].
- Analyze Boltzmann samples sets of structures [GASCOM10].
- Traverse the full mutational landscape [Recomb11].
- Explore enriched conformational spaces (PKs, MC-Fold) [...].
- Address constrained RNA design problems [...].



Using enumerative (models) and analytic (asymptotics) techniques, one can:

- Accurately enumerate RNA Shapes [JCB08].
- Improved random models for structures [Bioinf.06, TCS10].
- Analyze and optimize statistical sampling algorithms [JMB08].
- Analyze Boltzmann samples sets of structures [GASCOM10].
- Traverse the full mutational landscape [Recomb11].
- Explore enriched conformational spaces (PKs, MC-Fold) [...].
- Address constrained RNA design problems [...].





References I

-  F. Vanden Berghen.
CONDOR: a constrained, non-linear, derivative-free parallel optimizer for continuous, high computing load, noisy objective functions.
PhD thesis, IRIDIA, Universite Libre de Belgique, 2005.
-  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.
-  A. Denise, O. Roques, and M. Ternier.
Random generation of words of context-free languages according to the frequencies of letters.
In D. Gardy and A. Mokkadem, editors, *Mathematics and Computer Science: Algorithms, Trees, Combinatorics and probabilities*, Trends in Mathematics, pages 113–125. Birkhäuser, 2000.
-  A. Denise and P. Zimmermann.
Uniform random generation of decomposable structures using floating-point arithmetic.
Theor. Comput. Sci., 218(2):233–248, 1999.
-  P. Flajolet, P. Zimmermann, and B. Van Cutsem.
Calculus for the random generation of labelled combinatorial structures.
Theoretical Computer Science, 132:1–35, 1994.
-  J.S. McCaskill.
The equilibrium partition function and base pair binding probabilities for RNA secondary structure.
Biopolymers, 29:1105–1119, 1990.
-  R. Nussinov and A.B. Jacobson.
Fast algorithm for predicting the secondary structure of single-stranded RNA.
Proc Natl Acad Sci U S A, 77:6903–13, 1980.

References II



M. S. Waterman.

Secondary structure of single stranded nucleic acids.

Advances in Mathematics Supplementary Studies, 1(1):167–212, 1978.