

RNA as a combinatorial object

Asymptotics of RNA Shapes

Yann Ponty

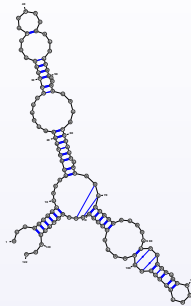
Bioinformatics Team
École Polytechnique/CNRS/INRIA AMIB – France

November 27, 2009

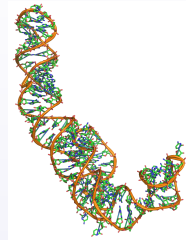
RNA structure

```
UUAGGCGGCCACAGC
GGUGGGGUUGCCUCC
CGUACCCAUCCGAA
CACGGAAGAUAGCC
CACCA GCGUCCGGG
GAGUACUGGAGUGCG
CGAGCCUCUGGGAAA
CCCGGUUCGCCGCCA
CC
```

Primary structure



Secondary structure



Tertiary structure

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

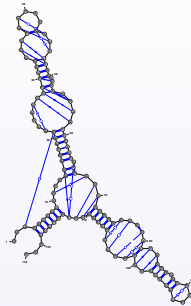
Definition

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

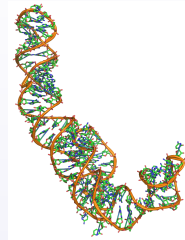
RNA structure

```
UUAGGCGGCCACAGC
GGUGGGGUUGCCUCC
CGUACCCAUCCGAA
CACGGAAGAUAGCC
CACCAGCGUCCGGG
GAGUACUGGAGUGCG
CGAGCCUCUGGGAAA
CCCGGUUCGCCGCCA
CC
```

Primary structure



Secondary⁺ structure



Tertiary structure

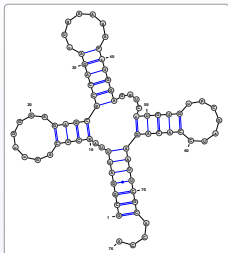
Source: 5s rRNA (PDBID: 1K73:B)

Definition

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

- 1 Foreword
 - Introduction
 - Motivation
- 2 Enumerative combinatorics 101
 - Generating functions
 - DSV/symbolic method
 - Singularity analysis
- 3 RNA shapes
 - Presentation
 - Motivation
 - π shapes

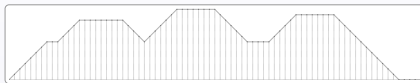
Various representations for a versatile molecule



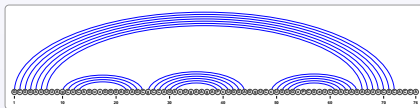
Outer planar graph

(((((((...(((.....))))((((.....))))))...((((.....))))))....

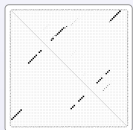
Well-parenthesized expression



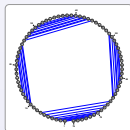
Mountain view



Non intersecting arcs



Dot plot

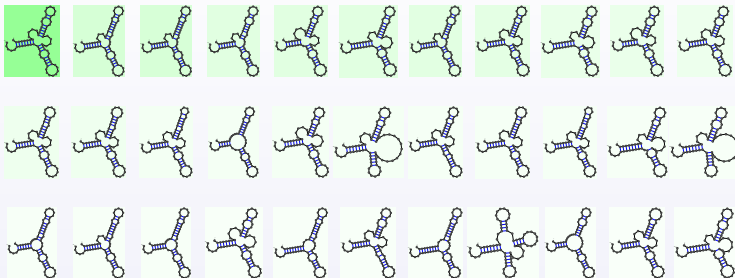


Feynman's diagram

Different objects
yet
Common combinatorial structure

Why use combinatorics?

Boltzmann ensemble is a (weighted) combinatorial class.



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

- Assess asymptotic properties of sec. str.
- Investigate worst and average-case complexities
- Obtain better algorithms for RNA

Generating functions

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

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

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

Main principles

Disclaimer

What follows, although true in this context, is embarassingly 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: $D(z) = \frac{1}{1-4z} \Rightarrow c_n \sim 4^n$
($\rho = 1/4, \alpha = -1, f(z) = 0$, 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  The diagram shows an equation between two wavy lines. The left side is a single wavy line with two dots at its ends. The right side is the sum of two terms: the first term is a wavy line with two dots at its ends and an additional dot on the line; the second term is a wavy line with two dots at its ends and a blue arc connecting two dots on the line. The equation is enclosed in a light blue box.

$$1 \cdot \text{---} = \text{---} \vee \text{---} \vee \epsilon$$

Appetizer: Motzkin words

Let us count dot-bracket notations (Motzkin words)

1	
2	$M \rightarrow \bullet M \mid (M)M \mid \epsilon$

Appetizer: Motzkin words

Let us count dot-bracket notations (Motzkin words)

1 

2 $M \rightarrow \bullet M \mid (M)M \mid \epsilon$

3
$$M(z) = z \cdot M(z) + z \cdot M(z) \cdot z \cdot M(z) + 1$$
$$= \frac{1 - z \pm \sqrt{1 - 2z - 3z^2}}{2z^2}$$

Appetizer: Motzkin words

Let us count dot-bracket notations (Motzkin words)

1 

2 $M \rightarrow \bullet M \mid (M)M \mid \epsilon$

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

Appetizer: Motzkin words

Let us count dot-bracket notations (Motzkin words)

$$1 \quad \text{---} = \text{---} \vee \text{---} \vee \epsilon$$

$$2 \quad M \rightarrow \bullet M \mid (M)M \mid \epsilon$$

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

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

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

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

RNA secondary structures

Consider RNA secondary structures (Waterman 78)

$$1 \quad \text{---} = \text{---} \vee \text{---} \overset{\geq 1}{\text{---}} \vee \epsilon$$

The diagram shows a recursive decomposition of RNA secondary structures. On the left, a wavy line with two dots is followed by an equals sign. To the right of the equals sign are three terms separated by 'or' symbols (\vee). The first term is a wavy line with two dots. The second term is a wavy line with two dots followed by a wavy line with two dots and a blue arc above it labeled with a greater-than-or-equal-to 1 (≥ 1). The third term is an empty set symbol (ϵ).

RNA secondary structures

Consider RNA secondary structures (Waterman 78)

1

$$\bullet \text{---} \bullet = \bullet \text{---} \bullet \text{---} \bullet \vee \bigvee_{\geq 1} \bullet \text{---} \bullet \vee \epsilon$$

2

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

RNA secondary structures

Consider RNA secondary structures (Waterman 78)

1

$$\bullet \text{---} \bullet = \bullet \text{---} \bullet \vee \text{---} \text{---} \vee \epsilon$$

≥ 1

2


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

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 

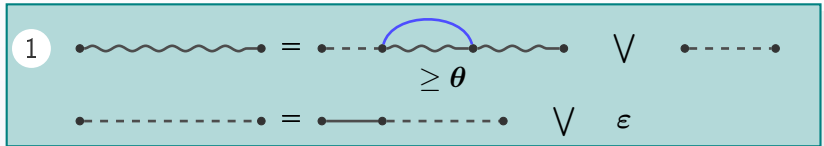
2
$$\begin{aligned} S &\rightarrow \bullet S \mid (\mathbf{T}) S \mid \epsilon \\ \mathbf{T} &\rightarrow \bullet S \mid (\mathbf{T}) S \end{aligned}$$

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

4
$$\rho = \frac{3 - \sqrt{5}}{2} = 1 - \phi$$
$$[z^n]S(z) = \sqrt{\frac{15 + 7\sqrt{5}}{8\pi}} \cdot \frac{\left(\frac{3 + \sqrt{5}}{2}\right)^n}{n\sqrt{n}} (1 + \mathcal{O}(1/n)) \sim 1.1 \cdot \frac{2.6^n}{n\sqrt{n}}$$

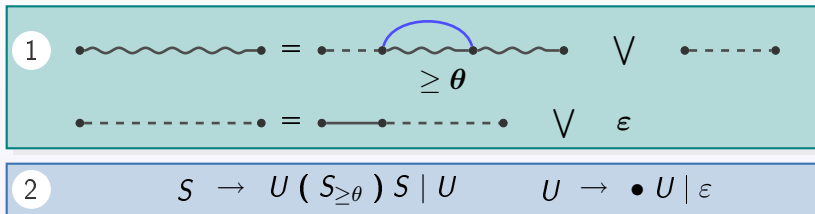
RNA secondary structures

Let us generalize the θ constraint



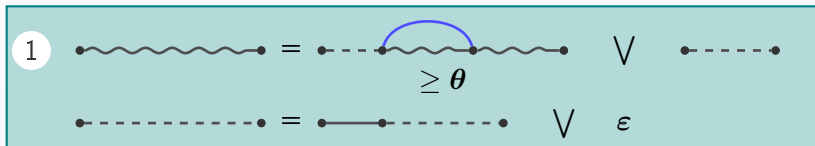
RNA secondary structures

Let us generalize the θ constraint



RNA secondary structures

Let us generalize the θ constraint

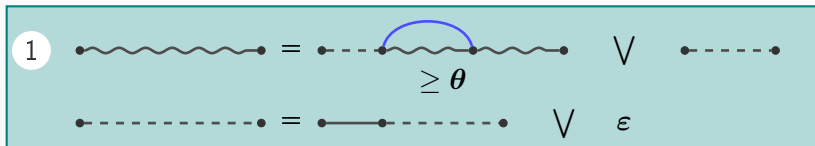


2

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

RNA secondary structures

Let us generalize the θ constraint



2

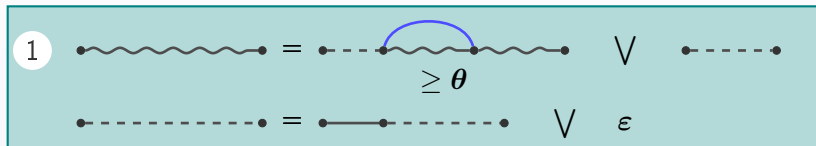
$$\begin{aligned} S &\rightarrow U(T)S \mid U & U &\rightarrow \bullet U \mid \epsilon \\ T &\rightarrow U(T)S \mid \bullet^\theta U \end{aligned}$$

3

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

RNA secondary structures

Let us generalize the θ constraint



2

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

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

Half-time report

Message #1

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

Message #2

Applying automatic theorems gives precise asymptotic equivalents.

Message #3

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

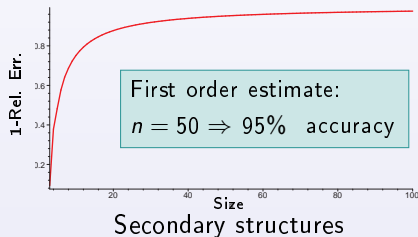
Half-time report

Message #1

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

Message #2

Applying **automatic theorems** gives **precise** asymptotic equivalents.



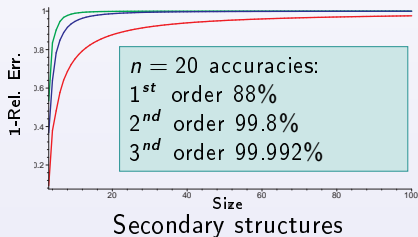
Half-time report

Message #1

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

Message #2

Applying **automatic theorems** gives **precise** asymptotic equivalents.



Half-time report

Message #1

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

Message #2

Applying **automatic theorems** gives **precise** asymptotic equivalents.

Message #3

There is a large exponential number of structures of size n :

Homopolymer model: $\Omega(2^n)$ **Stickiness model:** $\mathcal{O}(1.8^n/n^{3/2})$

- 1 Foreword
- 2 Enumerative combinatorics 101
- 3 RNA shapes
 - Presentation
 - Motivation
 - π shapes

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

Coarse-grain representation hierarchy for RNA sec. struct.

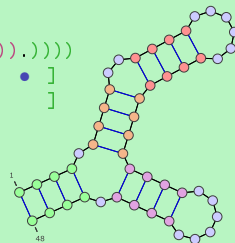
Based on the *underlying backbone* structure.

Example

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

π' -shape [• [• [•]] [•] •]

π -shape [[- -] []]



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

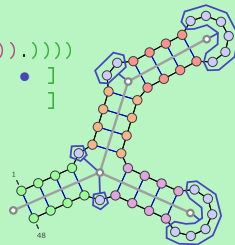
Coarse-grain representation hierarchy for RNA sec. struct.

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.

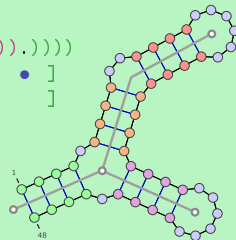
Example

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

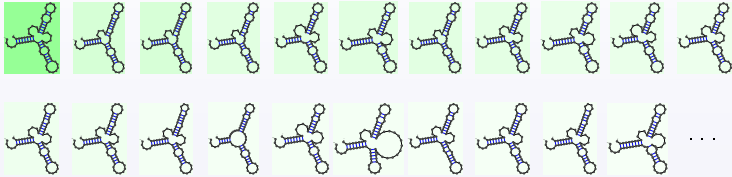
π' -shape [• [• [•]] [•] •]

π -shape [[- -] []]

Remove unpaired regions
Contract nested helices



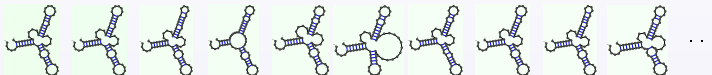
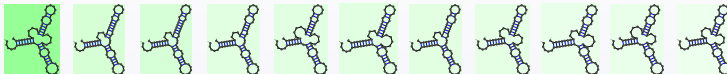
RNA shapes allow a hierarchical search in the Boltzmann ensemble



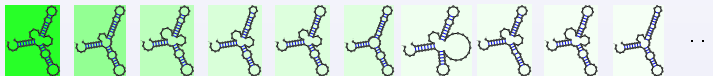
10000 samples \Rightarrow 1727 Secondary structures...

Motivation

RNA shapes allow a hierarchical search in the Boltzmann ensemble



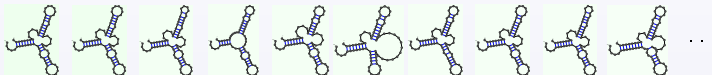
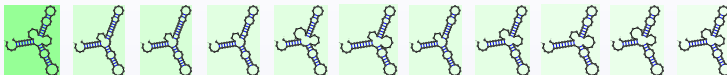
10000 samples \Rightarrow 1727 Secondary structures...



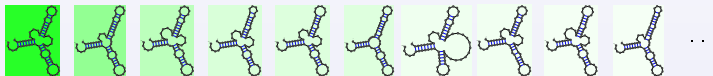
... 406 π' -shapes...

Motivation

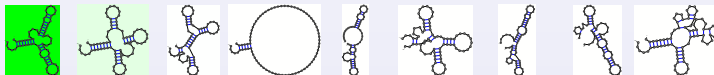
RNA shapes allow a hierarchical search in the Boltzmann ensemble



10000 samples \Rightarrow 1727 Secondary structures...



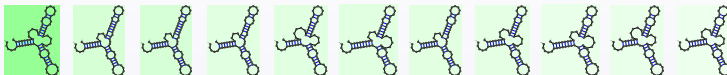
... 406 π' -shapes...



... but only 9 π -shapes!

Motivation

RNA shapes allow a hierarchical search in the Boltzmann ensemble

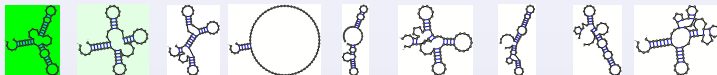


Is it reasonable to perform an exhaustive search of all possible shapes compatible with input structure?

10000 samples \Rightarrow 1727 Secondary structures

How many shapes must we investigate?

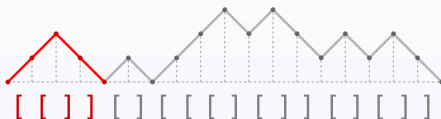
... 406 π' -shapes ...



... but only 9 π -shapes!

π -shapes

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

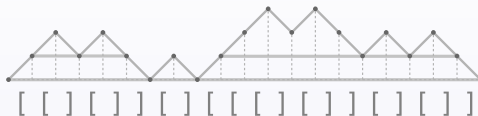


1

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

π -shapes

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

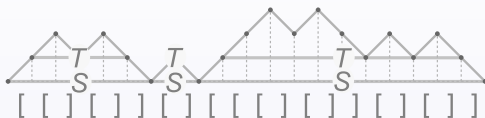


1 π -shapes are bracket words avoiding the $[[\dots]]$ motif.

2 $S \rightarrow [S/\{\dots\}]S \mid [S/\{\dots\}]$

π -shapes

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



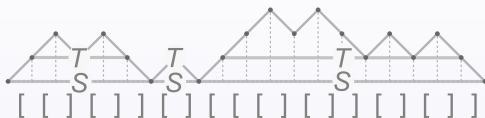
- 1 π -shapes are bracket words avoiding the $[[\dots]]$ motif.

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

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

π -shapes

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



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

2 $S \rightarrow [T]S \mid [T] \quad T \rightarrow [T]S \mid \varepsilon$

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???

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

$$2 \quad S \rightarrow [T]S \mid [T] \quad T \rightarrow [T]S \mid \varepsilon$$

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

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

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

$$\begin{aligned} 2 \quad S &\rightarrow [T]S \mid [T] & T &\rightarrow [T]S \mid \bullet^\theta \\ R &\rightarrow \square S \mid \varepsilon \end{aligned}$$

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

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

Limitations

Number of π -shapes of size n
 \neq
Number of π -shapes compatible with RNA of size n

Reasons:

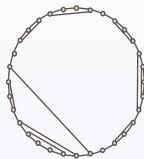
- 1 Shapes of size $\leq n$ should be considered
- 2 Forming a hairpin loop $[]$ takes at least $\theta + 2$ bases

$$\begin{aligned} 2 \quad S &\rightarrow [T]S \mid [T] & T &\rightarrow [T]S \mid \bullet^\theta \\ R &\rightarrow \square S \mid \varepsilon \end{aligned}$$

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

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

A surprising bijection



Theorem

$\#\pi$ shapes of size $n = \# \text{Motzkin words of length } 2n + 2$

Proof.

$$\begin{aligned} S(z) &= \frac{1-z^2-\sqrt{1-2z^2-3z^4}}{2z^2} & M(z) &= \frac{1-z-\sqrt{1-2z-3z^2}}{2z^2} \\ S(z) &= 1 + z^2 M(z^2) & \Rightarrow s_n &= m_{2n+2} \end{aligned}$$



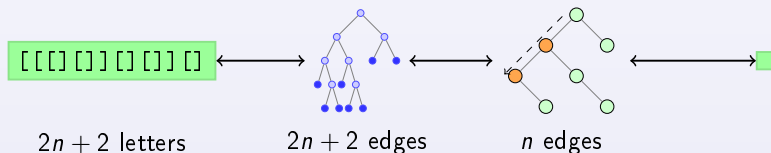
These two classes are in bijection.
How to state it? Can we exploit it?

Explicit bijection

Let $\psi, \phi : \{ [,] \}^* \rightarrow \{ (,) , \bullet \}$ such that

$$\begin{aligned}\psi((A) B) &= \begin{cases} \phi(A) & \text{If } B = \varepsilon \\ \phi(A) \bullet \psi(B) & \text{Otherwise} \end{cases} \\ \phi((A) B) &= \phi(A)[\psi(B)] \\ \phi(\varepsilon) &= \varepsilon.\end{aligned}$$

Then ψ is a **bijection** between s_{2n+2} and m_n .



Explicit bijection

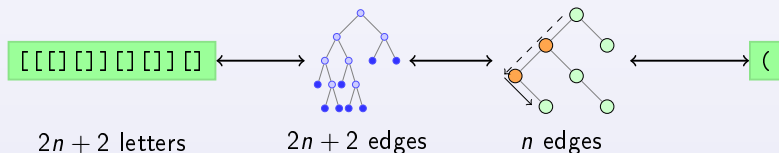
Let $\psi, \phi : \{ [,] \}^* \rightarrow \{ (,) , \bullet \}$ such that

$$\psi((A) B) = \begin{cases} \phi(A) & \text{If } B = \varepsilon \\ \phi(A) \bullet \psi(B) & \text{Otherwise} \end{cases}$$

$$\phi((A) B) = \phi(A)[\psi(B)]$$

$$\phi(\varepsilon) = \varepsilon.$$

Then ψ is a **bijection** between s_{2n+2} and m_n .

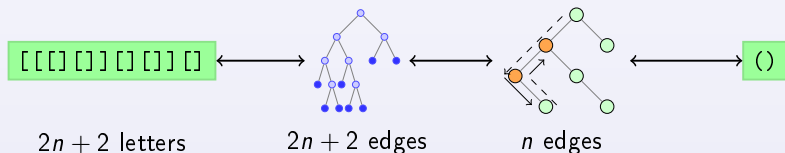


Explicit bijection

Let $\psi, \phi : \{ [,] \}^* \rightarrow \{ (,) , \bullet \}$ such that

$$\begin{aligned}\psi((A) B) &= \begin{cases} \phi(A) & \text{If } B = \varepsilon \\ \phi(A) \bullet \psi(B) & \text{Otherwise} \end{cases} \\ \phi((A) B) &= \phi(A)[\psi(B)] \\ \phi(\varepsilon) &= \varepsilon.\end{aligned}$$

Then ψ is a **bijection** between s_{2n+2} and m_n .



Explicit bijection

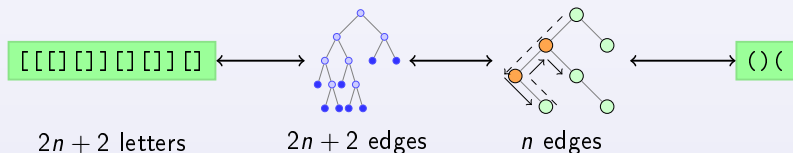
Let $\psi, \phi : \{ [,] \}^* \rightarrow \{ (,), \bullet \}$ such that

$$\psi((A) B) = \begin{cases} \phi(A) & \text{If } B = \varepsilon \\ \phi(A) \bullet \psi(B) & \text{Otherwise} \end{cases}$$

$$\phi((A) B) = \phi(A)[\psi(B)]$$

$$\phi(\varepsilon) = \varepsilon.$$

Then ψ is a **bijection** between s_{2n+2} and m_n .

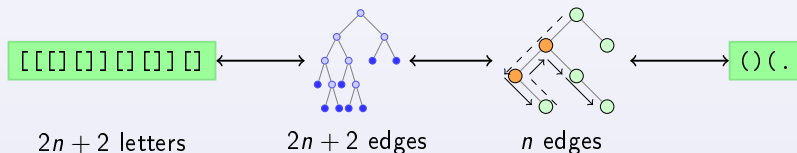


Explicit bijection

Let $\psi, \phi : \{ [,] \}^* \rightarrow \{ (,), \bullet \}$ such that

$$\begin{aligned}\psi((A) B) &= \begin{cases} \phi(A) & \text{If } B = \varepsilon \\ \phi(A) \bullet \psi(B) & \text{Otherwise} \end{cases} \\ \phi((A) B) &= \phi(A)[\psi(B)] \\ \phi(\varepsilon) &= \varepsilon.\end{aligned}$$

Then ψ is a **bijection** between s_{2n+2} and m_n .

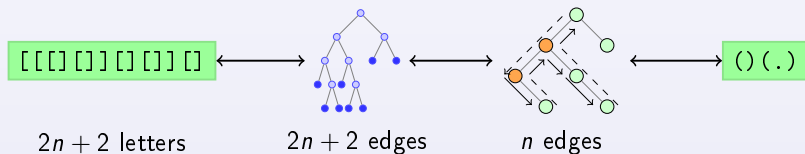


Explicit bijection

Let $\psi, \phi : \{ [,] \}^* \rightarrow \{ (,), \bullet \}$ such that

$$\begin{aligned}\psi((A) B) &= \begin{cases} \phi(A) & \text{If } B = \varepsilon \\ \phi(A) \bullet \psi(B) & \text{Otherwise} \end{cases} \\ \phi((A) B) &= \phi(A)[\psi(B)] \\ \phi(\varepsilon) &= \varepsilon.\end{aligned}$$

Then ψ is a **bijection** between s_{2n+2} and m_n .



Explicit bijection

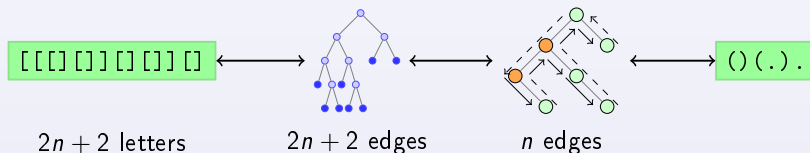
Let $\psi, \phi : \{ [,] \}^* \rightarrow \{ (,) , \bullet \}$ such that

$$\psi((A) B) = \begin{cases} \phi(A) & \text{If } B = \varepsilon \\ \phi(A) \bullet \psi(B) & \text{Otherwise} \end{cases}$$

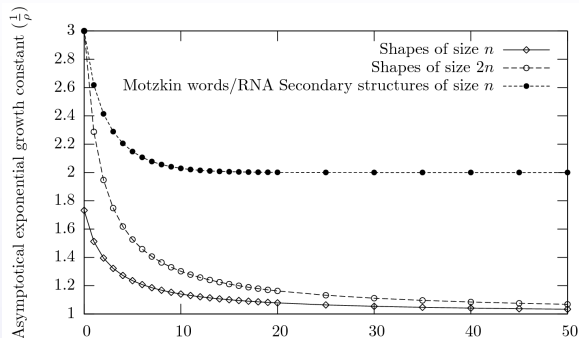
$$\phi((A) B) = \phi(A)[\psi(B)]$$

$$\phi(\varepsilon) = \varepsilon.$$

Then ψ is a **bijection** between s_{2n+2} and m_n .



Limits of the bijection



Impacts of θ on shapes and Motzkin are drastically different.

Theorem

Expectations of number of term. loops in Motzkin words and π -shapes scale like $m_n^t \sim \frac{n}{6} + \mathcal{O}(1)$ and $s_{2n+2}^t \sim \frac{2n}{3} + \mathcal{O}(1)$

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

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

2
$$\begin{aligned} R &\rightarrow \square R | S & S &\rightarrow U [T] S | U & U &\rightarrow \diamond | \varepsilon \\ T &\rightarrow U [T] U [T] S | \diamond [T] | [T] \diamond | \diamond [T] \diamond | \bullet^\theta \end{aligned}$$

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
$$\begin{aligned} R &\rightarrow \square R \mid S & S &\rightarrow U[T]S \mid U & U &\rightarrow \diamond \mid \varepsilon \\ T &\rightarrow U[T]U[T]S \mid \diamond[T] \mid [T]\diamond \mid \diamond[T]\diamond \mid \bullet^\theta \end{aligned}$$

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

Summary

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	$1.28 \cdot \frac{1.81^n}{n\sqrt{n}}$

- For *context-free* objects, finding gen. fun. is **easy**...
...and **precise asymptotics estimates** follow readily
- **Bijection** between Motkzin words and π -shapes
- **Way less** many shapes than sec. str.!
- Homopolymer model **overestimates** number of shapes
Need for a probabilistic model for base-pairing
But **stickiness** is not enough...

Collaborators: W. A. Lorenz and P. Clote (Boston College)