# Sampling
# RNA Secondary Structures
# with Pseudoknots
# using Analytic Combinatorics
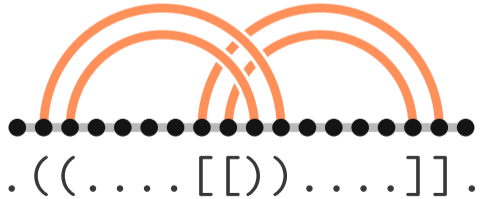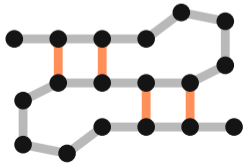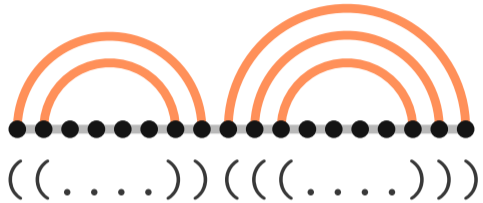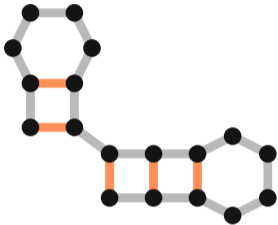
**Casper Asbjørn Eriksen[1]**  **Daniel Merkle[1,2]**
**Markus Nebel[2]**  **Jonas Vistrup[1]**

[1]Algorithmic Cheminformatics Group, University of Southern Denmark
[2]Faculty of Technology, Bielefeld University

**39th TBI Winterseminar**

# RNA pseudoknots

```
((....))(((....)))
```

```
.((....[[))....]].
```

SDU

# Grammars for RNA Secondary Structures

The following CFG generates well-formed dot-bracket strings (Motzkin words):

$$S \rightarrow \bullet S \mid (S) S \mid \epsilon$$

However, we might want a grammar with biologically meaningful rules.

- Energy estimations
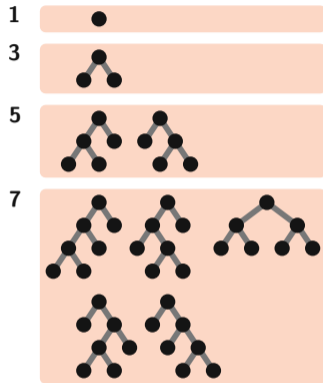- Parametrised generation
- *Uniform / non-uniform Boltzmann sampling*

$$
\begin{aligned}
&f_1 = S \rightarrow TAC, \\
&f_2 = T \rightarrow TAC, \\
&f_3 = T \rightarrow C, \\
&f_4 = C \rightarrow C \bullet, \\
&f_5 = C \rightarrow \epsilon,
\end{aligned}
\right\} \text{ exterior loop}
$$

$$
\begin{aligned}
&f_6 = A \rightarrow (L), \\
&f_7 = L \rightarrow (L),
\end{aligned}
\right\} \text{ initiate and extend stem}
$$

$f_8 = L \rightarrow M, \quad \text{initiate multiloop}$

$$
\begin{aligned}
&f_9 = L \rightarrow P, \\
&f_{10} = L \rightarrow Q,
\end{aligned}
\right\} \text{ initiate interior loop}
$$

$f_{11} = L \rightarrow R,$

$f_{12} = L \rightarrow F, \quad \text{initiate hairpin loop}$

$f_{13} = L \rightarrow G, \quad \text{initiate bulge loop}$

$$
\begin{aligned}
&f_{14} = G \rightarrow (L) \bullet, \\
&f_{15} = G \rightarrow (L) B \bullet \bullet, \\
&f_{16} = G \rightarrow \bullet (L), \\
&f_{17} = G \rightarrow \bullet \bullet B(L), \\
&f_{18} = B \rightarrow B \bullet, \\
&f_{19} = B \rightarrow \epsilon,
\end{aligned}
\right\} \text{ bulge loops}
$$

$$
\begin{aligned}
&f_{20} = F \rightarrow \bullet \bullet \bullet, \\
&f_{21} = F \rightarrow \bullet \bullet \bullet \bullet, \\
&f_{22} = F \rightarrow \bullet \bullet \bullet \bullet \bullet H, \\
&f_{23} = H \rightarrow H \bullet, \\
&f_{24} = H \rightarrow \epsilon,
\end{aligned}
\right\} \text{ hairpin loop}
$$

$$
\begin{aligned}
&f_{25} = P \rightarrow \bullet (L) \bullet, \\
&f_{26} = P \rightarrow \bullet (L) \bullet \bullet, \\
&f_{27} = P \rightarrow \bullet \bullet (L) \bullet, \\
&f_{28} = P \rightarrow \bullet \bullet (L) \bullet \bullet,
\end{aligned}
\right\} \text{ small interior loops}
$$

$$
\begin{aligned}
&f_{29} = Q \rightarrow \bullet \bullet (L) K \bullet \bullet \bullet, \\
&f_{30} = Q \rightarrow \bullet \bullet \bullet J(L) K \bullet \bullet, \\
&f_{31} = R \rightarrow \bullet (L) K \bullet \bullet \bullet, \\
&f_{32} = R \rightarrow \bullet \bullet \bullet J(L) \bullet, \\
&f_{33} = J \rightarrow J \bullet, \\
&f_{34} = J \rightarrow \epsilon, \\
&f_{35} = K \rightarrow K \bullet, \\
&f_{36} = K \rightarrow \epsilon,
\end{aligned}
\right\} \text{ other interior loops}
$$

$$
\begin{aligned}
&f_{37} = M \rightarrow U(L)U(L)N, \\
&f_{38} = N \rightarrow U(L)N, \\
&f_{39} = N \rightarrow U, \\
&f_{40} = U \rightarrow U \bullet, \\
&f_{41} = U \rightarrow \epsilon.
\end{aligned}
\right\} \text{ multiloop}
$$

(Nebel et. al 2011)

# Analytic Combinatorics

(Ordered) Binary trees:

$$B = \bullet + [\bullet \times B \times B]$$



*"If you can specify it, you can analyse it."*
*- Philippe Flajolet*

SDU

(Ordered) Binary trees:

$$B = \bullet + [\bullet \times B \times B]$$

Number of binary trees of size $n$:

$1, 0, 1, 0, 2, 0, 5, 0, 14\ldots$   (seq. **A000108**)

**1**


**3**


**5**


**7**


*"If you can specify it, you can analyse it."*
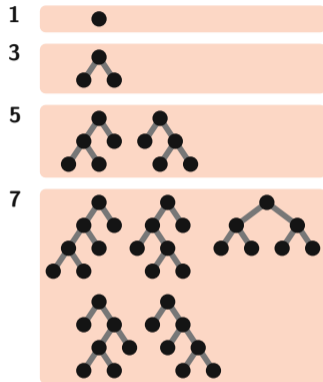*- Philippe Flajolet*

# Analytic Combinatorics

(Ordered) Binary trees:

$$B = \bullet + [\bullet \times B \times B]$$
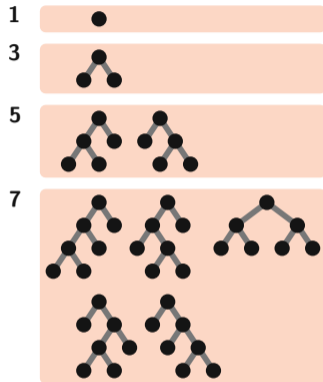
Number of binary trees of size $n$:

$1, 0, 1, 0, 2, 0, 5, 0, 14...$ (seq. **A000108**)

**Generating function:**

$GF_B(z) = z + zGF_B(z)^2$   Symbolic transfer thm.

$\quad = \dfrac{1 - \sqrt{1 - 4z}}{2}$   Solve + expansion

$\quad = \mathbf{1}z + \mathbf{1}z^3 + \mathbf{2}z^5 + \mathbf{5}z^7 + \mathbf{14}z^9 + \cdots$



*"If you can specify it, you can analyse it."*
*- Philippe Flajolet*

SDU

Description of combinatorial structures:

- ◆ Binary trees: $B = \bullet + [\bullet \times B \times B]$
- ◆ General trees: $T = \bullet + [\bullet \times \text{SEQ}(T)]$
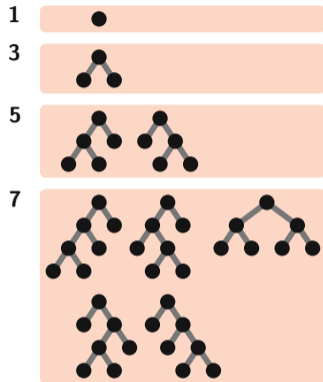- ◆ Derangements: $D = \text{SET}(\text{CYC}_{\geq 1}(\bullet))$



*"If you can specify it, you can analyse it."*
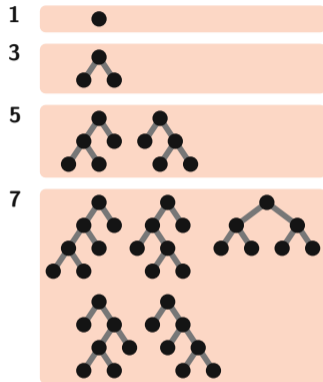*- Philippe Flajolet*

SDU

# Analytic Combinatorics

Description of combinatorial structures:

- Binary trees: $B = \bullet + [\bullet \times B \times B]$
- General trees: $T = \bullet + [\bullet \times \text{SEQ}(T)]$
- Derangements: $D = \text{SET}(\text{CYC}_{\geq 1}(\bullet))$

A context-free language is a combinatorial class:

- Grammar: $S \rightarrow \bullet S \mid (S)\,S \mid \epsilon$
- Specification: $S = [\bullet \times S] + [\,(\times S \times\,)\,\times S]$



*"If you can specify it, you can analyse it."*
*- Philippe Flajolet*

SDU

# Analytic Combinatorics

A **Boltzmann Sampler** samples class a uniformly

- ◆ Recursive algorithm using GF
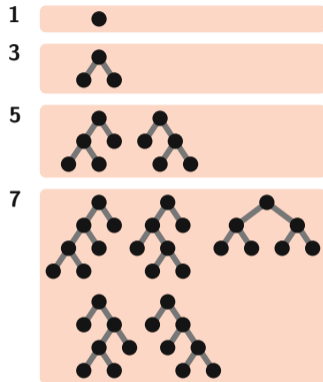


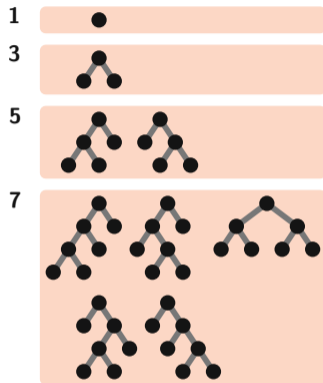*"If you can specify it, you can analyse it."*
*- Philippe Flajolet*

SDU

# Analytic Combinatorics

A **Boltzmann Sampler** samples class a uniformly

- ◆ Recursive algorithm using GF

**Maximum Likelihood Sampling**

- ◆ Non-uniform sampling by weighing construction rules
- ◆ Obtain weights by parsing an ensemble of structures



*"If you can specify it, you can analyse it."*
*- Philippe Flajolet*

SDU

# Boltzmann Combinatorial Object Sampler

**BCOS** is a *work-in progress* Boltzmann sampling libary for:

- Uniform sampling of arbitrary combinatorial classes
- Weighted sampling with predefined weights
- Maximum Likelihood ensemble-based sampling

SDU

**BCOS** is a *work-in progress* Boltzmann sampling libary for:

- Uniform sampling of arbitrary combinatorial classes
- Weighted sampling with predefined weights
- Maximum Likelihood ensemble-based sampling

**Features**:

- Dynamic arbitrary-precision evaluation of generating functions guarantees correctness
- Ability to sample complex classes even if no closed form GF can be found.
- Python interface for Combinatorial object semantics allows intuitive integration.

**Ordered tree**

```python
import bcos
tree_class = bcos.System(
    "T = N + (N * SEQ(T)),
     N = Atom"
    )
tree_class.sample(size=(5,10))

> (N,((N,(N,N,N,N)),(N,(N))))
```

## Ordered tree

```
import bcos
tree_class = bcos.System(
    "T = N + (N * SEQ(T)),
     N = Atom"
    )
tree_class.sample(size=(5,10))

> (N,((N,(N,N,N,N)),(N,(N))))
```

## RNA secondary structure

```
import bcos
g = bcos.cfg(
    "S -> . S | ( S ) | "
    )
gclass = bocs.System(g)
gclass.sample(size=21)

> ..((.(..)))..((.))(.)
```

# Multiple Context-Free Grammars

However, CFGs are not expressive enough to describe secondary structures with pseudoknots.

Introducing: **Multiple Context Free Grammars** (MCFG)

- Describes a broader class of languages than CFGs
- Allows for non-local correlations

Next step: designing an MCFG for pseudoknot structures.

Example MCFG
( $a^n b^n c^n d^n$ ):

$$\begin{pmatrix} A_1 \\ A_2 \end{pmatrix} \rightarrow \begin{pmatrix} ab \\ cd \end{pmatrix}$$

$$\begin{pmatrix} A_1 \\ A_2 \end{pmatrix} \rightarrow \begin{pmatrix} aA_1b \\ cA_2d \end{pmatrix}$$

Example derivation:

$$A_1A_2$$
$$\rightarrow aA_1bcA_2d$$
$$\rightarrow aabbccdd$$

SDU

# Pseudoknot Grammars

Problem: Potentially infinite alphabet

- $(), [], \{\}, <>, Aa, Bb, \cdots \quad \cdots, \mathrm{A}\alpha, \mathrm{B}\beta, \Gamma\gamma$?

Problem: Potentially infinite alphabet

- $(), [], \{\}, <>, Aa, Bb, \cdots \quad \cdots, \mathrm{A}\alpha, \mathrm{B}\beta, \Gamma\gamma?$

Possible solution:

- Parametrised alphabet:

$$S \to A_1 A_2 \mid \bullet \mid (_X\, S\, )_X$$

$$\begin{pmatrix} A_1 \\ A_2 \end{pmatrix} \to \begin{pmatrix} SA_1 \\ A_2 \end{pmatrix} \,\Big|\, \begin{pmatrix} A_1 \\ SA_2 \end{pmatrix} \,\Big|\, \begin{pmatrix} (_X A_1 \\ )_X A_2 \end{pmatrix} \,\Big|\, \begin{pmatrix} \epsilon \\ \epsilon \end{pmatrix}$$
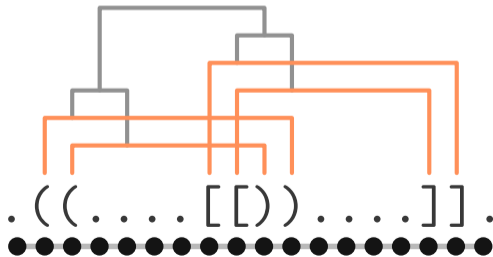
where $(_X, )_X = (), [], \cdots.$

# Pseudoknot Grammars

Parsing:
- Instantiate grammar with sufficient degree and parser
  **or**:
- Generate the lowest level of the parse tree and parse this

Generation:
- Generate using simple grammar and create string from parse tree

# Sampling Pseudoknot Structures

**Strategy**

- Let $G$ be a parametrised MCFG (**PMCFG**)
- Let $L$ be a library of pseudoknot RNA dot-bracket strings, possibly containing pseudoknots

SDU

### Strategy

- Let *G* be a parametrised MCFG (**PMCFG**)
- Let *L* be a library of pseudoknot RNA dot-bracket strings, possibly containing pseudoknots
- Instantiate *G* and parse each string in *L*

# Sampling Pseudoknot Structures

### Strategy

- Let *G* be a parametrised MCFG (**PMCFG**)
- Let *L* be a library of pseudoknot RNA dot-bracket strings, possibly containing pseudoknots
- Instantiate *G* and parse each string in *L*
- Result is a stochastic parametrised MCFG (**SPMCFG**)

SDU

### Strategy

- Let *G* be a parametrised MCFG (**PMCFG**)
- Let *L* be a library of pseudoknot RNA dot-bracket strings, possibly containing pseudoknots
- Instantiate *G* and parse each string in *L*
- Result is a stochastic parametrised MCFG (**SPMCFG**)
- Convert to a weighted combinatorial class *R*

# Sampling Pseudoknot Structures

### Strategy

- Let $G$ be a parametrised MCFG (**PMCFG**)
- Let $L$ be a library of pseudoknot RNA dot-bracket strings, possibly containing pseudoknots
- Instantiate $G$ and parse each string in $L$
- Result is a stochastic parametrised MCFG (**SPMCFG**)
- Convert to a weighted combinatorial class $R$
- Calculate GF based on **SCFG** reduction

# Sampling Pseudoknot Structures

### Strategy

- Let $G$ be a parametrised MCFG (**PMCFG**)
- Let $L$ be a library of pseudoknot RNA dot-bracket strings, possibly containing pseudoknots
- Instantiate $G$ and parse each string in $L$
- Result is a stochastic parametrised MCFG (**SPMCFG**)
- Convert to a weighted combinatorial class $R$
- Calculate GF based on **SCFG** reduction
- Sample $R$, and create dot-bracket string based on the parse tree.

SDU