# Combinatorics of Biomolecules

## C.M. Reidys

### Nankai University
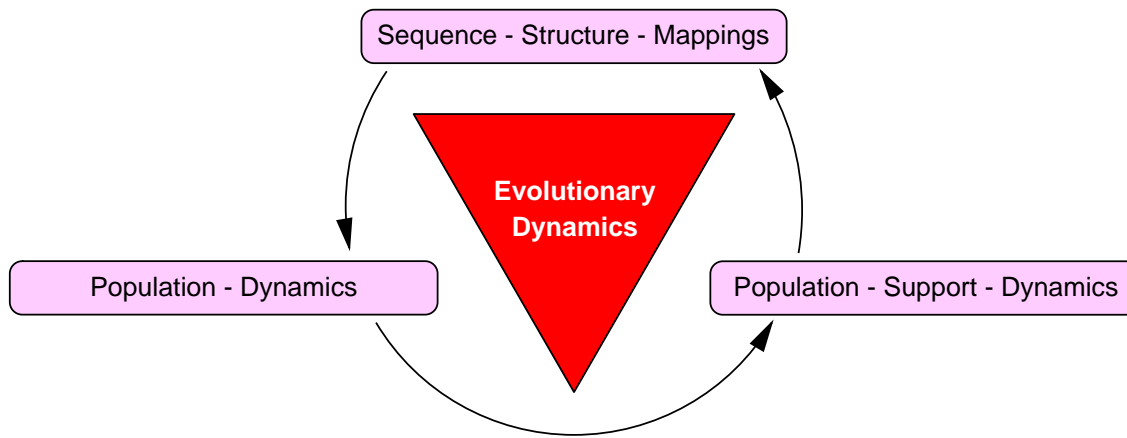
### Center for Combinatorics, LMPC

FIGURE 1. Evolutionary Dynamics

# Computational Biology Group at Nankai

- sequence to structure maps
- combinatorial representation of biomolecules
- new generation folding algorithms of biomolecules
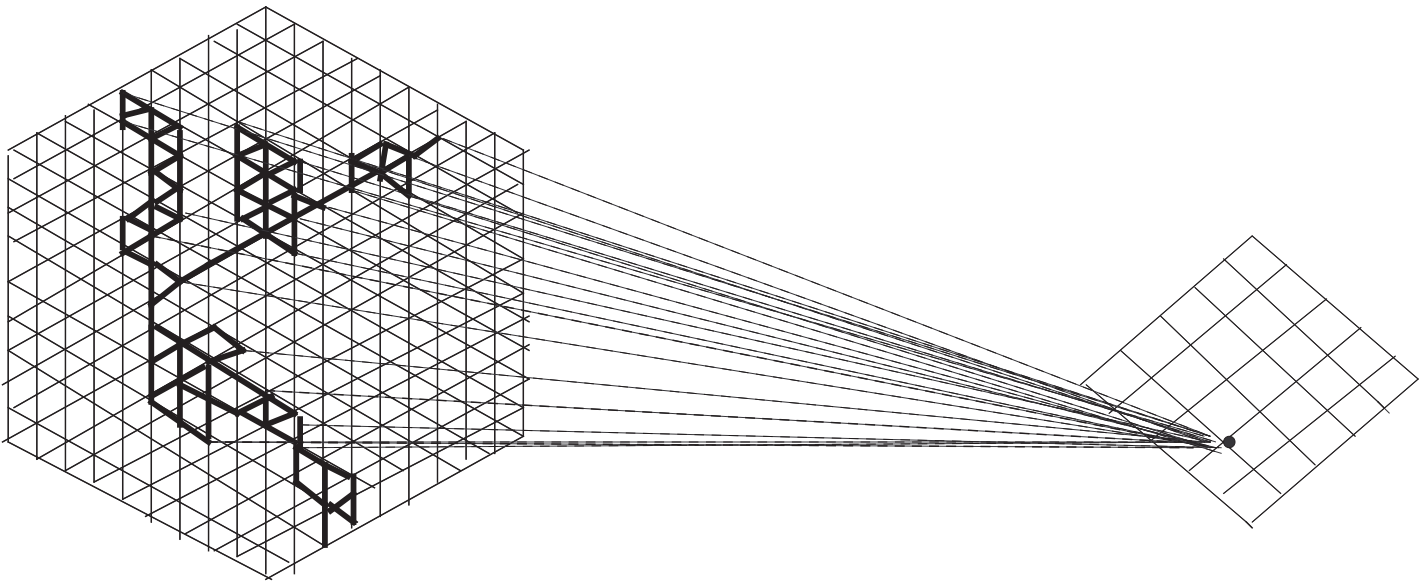
# Sequences and Shapes



FIGURE 2. The neutral network of a structure. Sequence space (right) and shape space (left) represented as lattices. We draw the edges between two sequences bold if they map into the one particular structure on the left. The two key properties of neutral nets are their connectivity and percolation. They allow sequences to move while maintaining a shape through sequence space.
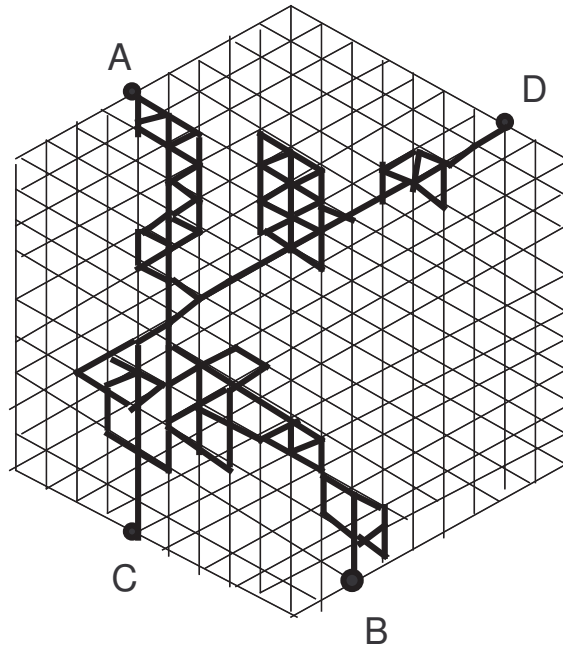
# Sequences and Shapes: Neutral Networks



FIGURE 3. Neutral network. Sequence space is represented as lattice and the neutral net is an induced subgraph (bold edges). We label the pairs of sequences representing antipodal pairs by $(A, B)$ and $(C, D)$. The two key properties of neutral nets are their connectivity and percolation.

**Theorem 1.** *Let $Q_{2,\lambda_n}^n$ be the random graph consisting of $Q_2^n$-subgraphs, $\Gamma_n$, induced by selecting each $Q_2^n$-vertex with independent probability $\lambda_n = \frac{1+\chi_n}{n}$, where $\chi_n = \epsilon n^{\frac{a-1}{2}}$, where $0 < \epsilon$ and $0 < a \leq 1$. Then we have*

$$(0.1) \qquad \exists\, \kappa_a > 0; \quad \lim_{n \to \infty} \mathbb{P}\left( |C_n^{(1)}| \geq \kappa_a\, n^{a-1} |\Gamma_n| \text{ and } C_n^{(1)} \text{ is unique} \right) = 1 .$$

Christian M. Reidys *Large components in random induced subgraphs of n-cubes* Discrete Math. submitted, 2007.
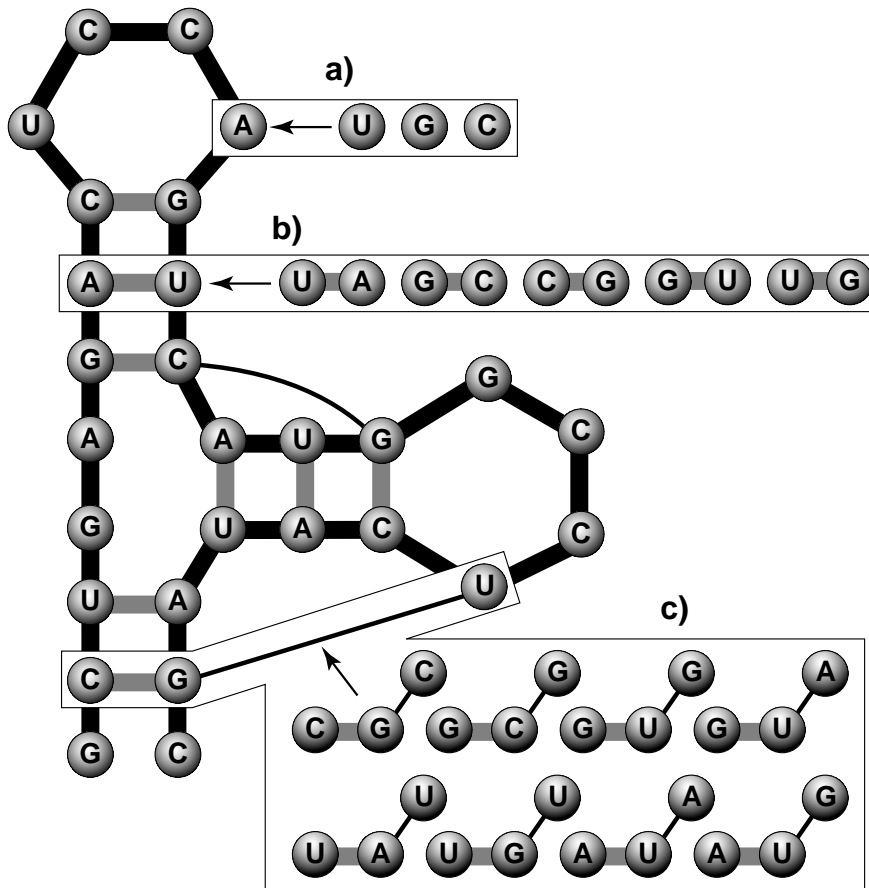
FIGURE 4. RNA secondary structure. Watson-Crick base-pairs (gray), tertiary contacts (black)

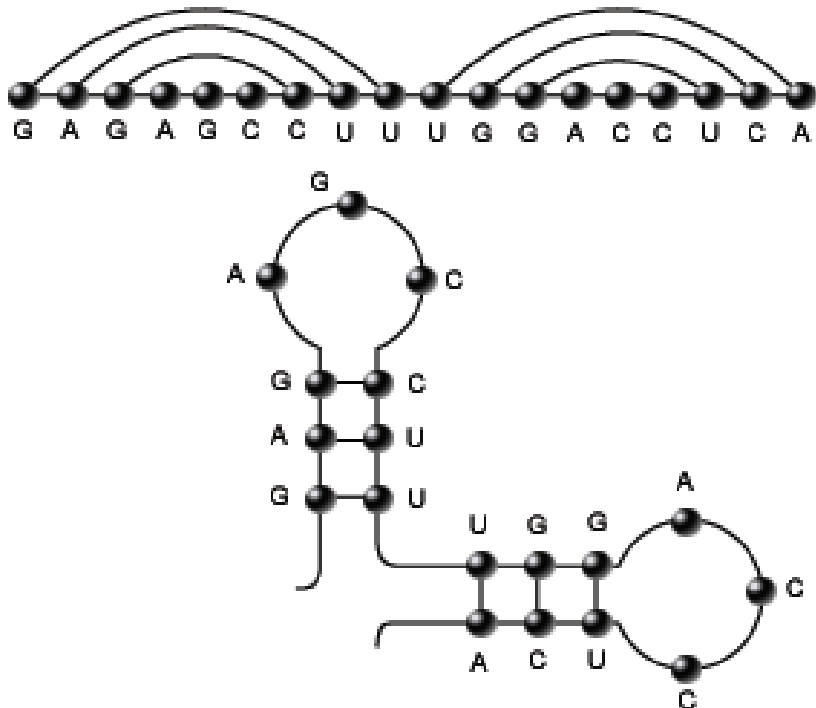# RNA secondary structures or better: 2-noncrossing RNA



FIGURE 5. RNA secondary structures. Diagram representation (top): the primary sequence, **GAGAGCCUUUGGACCUCA**, is drawn horizontally and its backbone bonds are ignored. All bonds are drawn in the upper halfplane and secondary structures have the property that no two arcs intersect and all arcs have minimum length 2. Outer planar graph representation (bottom).
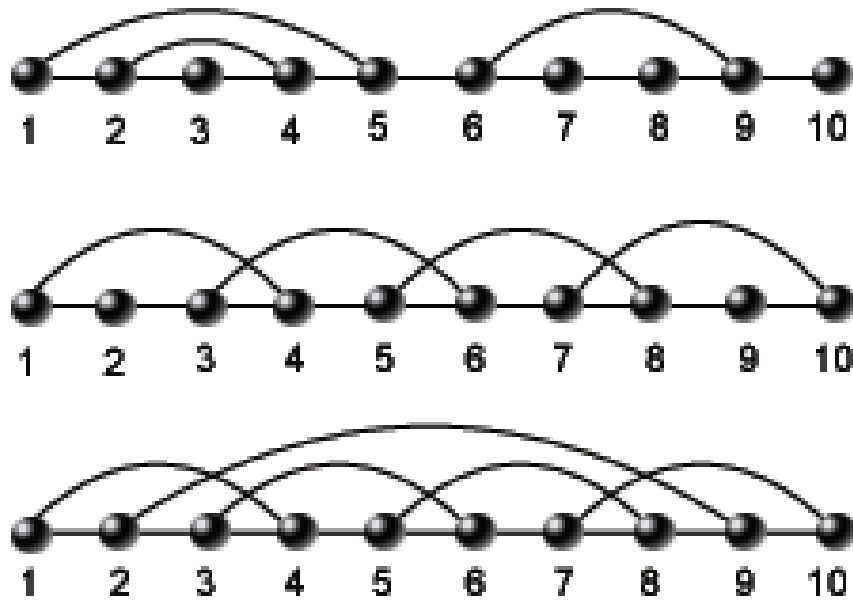
# 3-noncrossing RNA structures



FIGURE 6. $k$-noncrossing RNA structures. (a) **secondary structure**, (b) **planar** 3-**noncrossing RNA structure**, (c) the smallest **non-planar** 3-**noncrossing structure**

**Definition 1.** An RNA structure (of pseudoknot type $k - 2$), $S_{k,n}$, is a digraph in which all vertices have degree $\leq 1$, that does not contain a $k$-set of mutually intersecting arcs and 1-arcs, i.e. arcs of the form $(i, i + 1)$, respectively.

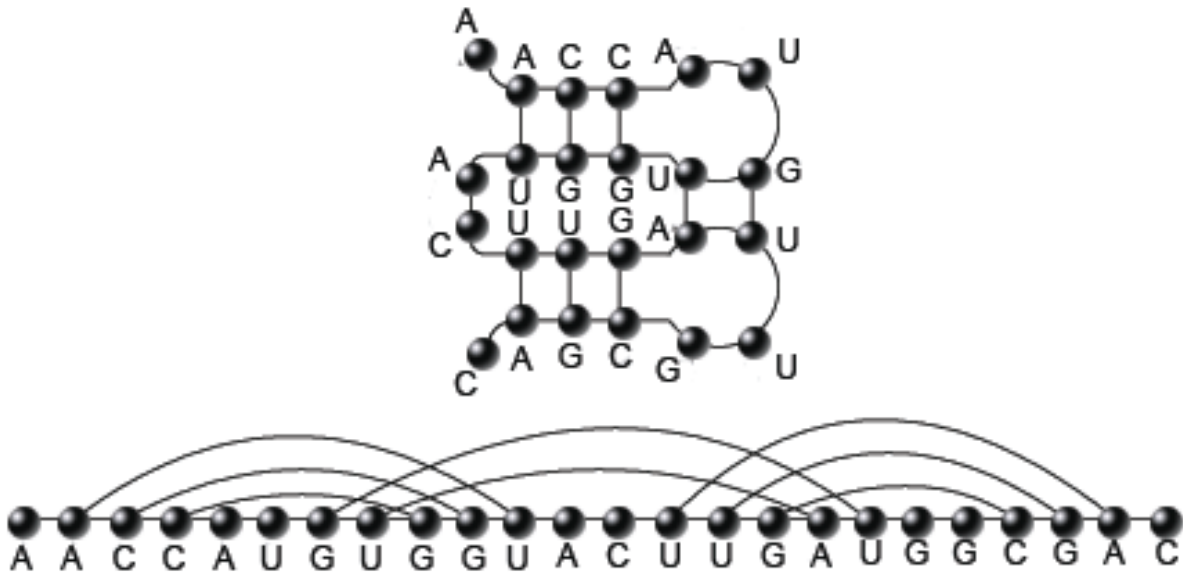# 3-noncrossing RNA structures: What is new?



FIGURE 7. A 3-noncrossing RNA structure, as a planar graph (top) and as a diagram (bottom)



FIGURE 8. The proposed SRV-1 frame-shift is a 10-noncrossing RNA structure motif.

# Combinatorics of $3$-noncrossing RNA structures

**Theorem 2.** *Let $k \in \mathbb{N}$, $k \geq 2$, let $f_k(n, \ell)$ be the number of $k$-noncrossing digraphs over $n$ vertices with exactly $\ell$ isolated vertices. Then the number of RNA structures with $\ell$ isolated vertices, $S_k(n, \ell)$, is*

$$(0.2) \qquad S_k(n, \ell) = \sum_{b=0}^{(n-\ell)/2} (-1)^b \binom{n-b}{b} f_k(n - 2b, \ell) .$$

*Furthermore the number of $k$-noncrossing RNA structures, $S_k(n)$ is given by*

$$(0.3) \qquad S_k(n) = \sum_{b=0}^{\lfloor n/2 \rfloor} (-1)^b \binom{n-b}{b} \left\{ \sum_{\ell=0}^{n-2b} f_k(n - 2b, \ell) \right\}$$

Emma Y. Jin, Jing Qin and Christian M. Reidys *Combinatorics of RNA Structures with Pseudoknots*, Bulletin of Math. Bio., 2007, in press.

| $n$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $S_3(n)$ | 1 | 1 | 2 | 5 | 13 | 36 | 105 | 321 | 1018 | 3334 | 11216 | 38635 | 135835 | 486337 | 1769500 |

**Table 1.** The first $15$ numbers of $3$-noncrossing RNA structures.

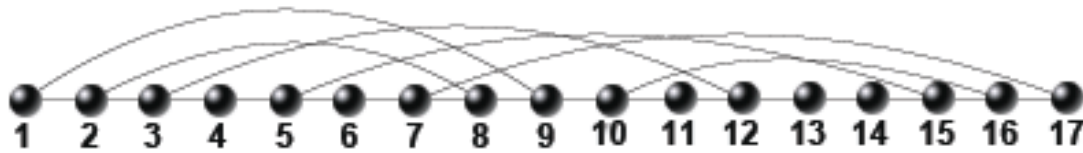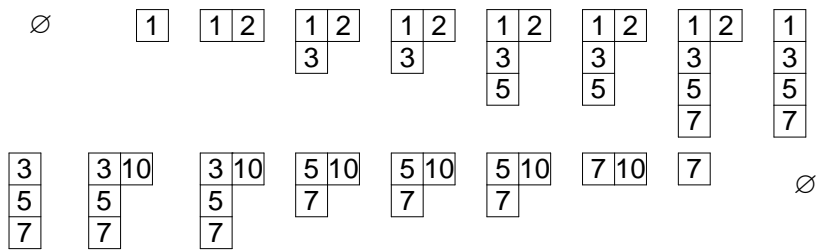# Combinatorics of 3-noncrossing RNA structures: Main idea



FIGURE 9. A 5-noncrossing structure corresponding to the oscillating tableau below and subsequently the corresponding walk $\gamma_{a,\bar{a}}$ in $\mathbb{Z}^4$.

# Why 3-noncrossing RNA structures is so different: recursions

**Corollary 1.** *The number of RNA secondary structures having exactly $\ell$ isolated vertices,* $S_2(n, \ell)$, *is given by*

(0.4)
$$S_2(n, \ell) = \frac{2}{n - \ell} \binom{\frac{n+\ell}{2}}{\frac{n-\ell}{2} + 1} \binom{\frac{n+\ell}{2} - 1}{\frac{n-\ell}{2} - 1}.$$

*Furthermore* $S_2(n, \ell)$ *satisfies the recursion*

(0.5)
$$(n - \ell)(n - \ell + 2) \cdot S_2(n, \ell) - (n + \ell)(n + \ell - 2) \cdot S_2(n - 2, \ell) = 0.$$

**Corollary 2.** *The number of* 3-*noncrossing RNA structures having exactly $\ell$ isolated vertices,* $S_3(n, \ell)$, *satisfies the* 4-*term recursion*

(0.6)
$$p_1(n, \ell)\, S_3(n - 6, \ell) - p_2(n, \ell)\, S_3(n - 4, \ell) - p_3(n, \ell) S_3(n - 2, \ell) + p_4(n, \ell)\, S_3(n, \ell) = 0,$$

*where the coefficients* $p_1(n, \ell)$, $p_2(n, \ell)$ $p_3(n, \ell)$ *and* $p_4(n, \ell)$ *are given by*

$$p_1(n, \ell) = \frac{1}{2}n(n - 1)(n - 10 + \ell)(n - 4 + \ell)(n - 8 + \ell)$$

$$p_2(n, \ell) = \frac{1}{2}n(n - 3)(13n^3 - 126n^2 + 13n^2\ell - 88n\ell + 392n + 3n\ell^2 + 216\ell - 384 - 42\ell^2 + 3\ell^3)$$

$$p_3(n, \ell) = (n - 1)(\frac{1}{2}n - 2)(13n^3 - 30n^2 - 13n^2\ell + 8n + 16n\ell + 3n\ell^2 + 30\ell^2 - 72\ell - 3\ell^3)$$

$$p_4(n, \ell) = (n - 3)(\frac{1}{2}n - 2)(n - \ell)(n - \ell + 6)(n - \ell + 4).$$

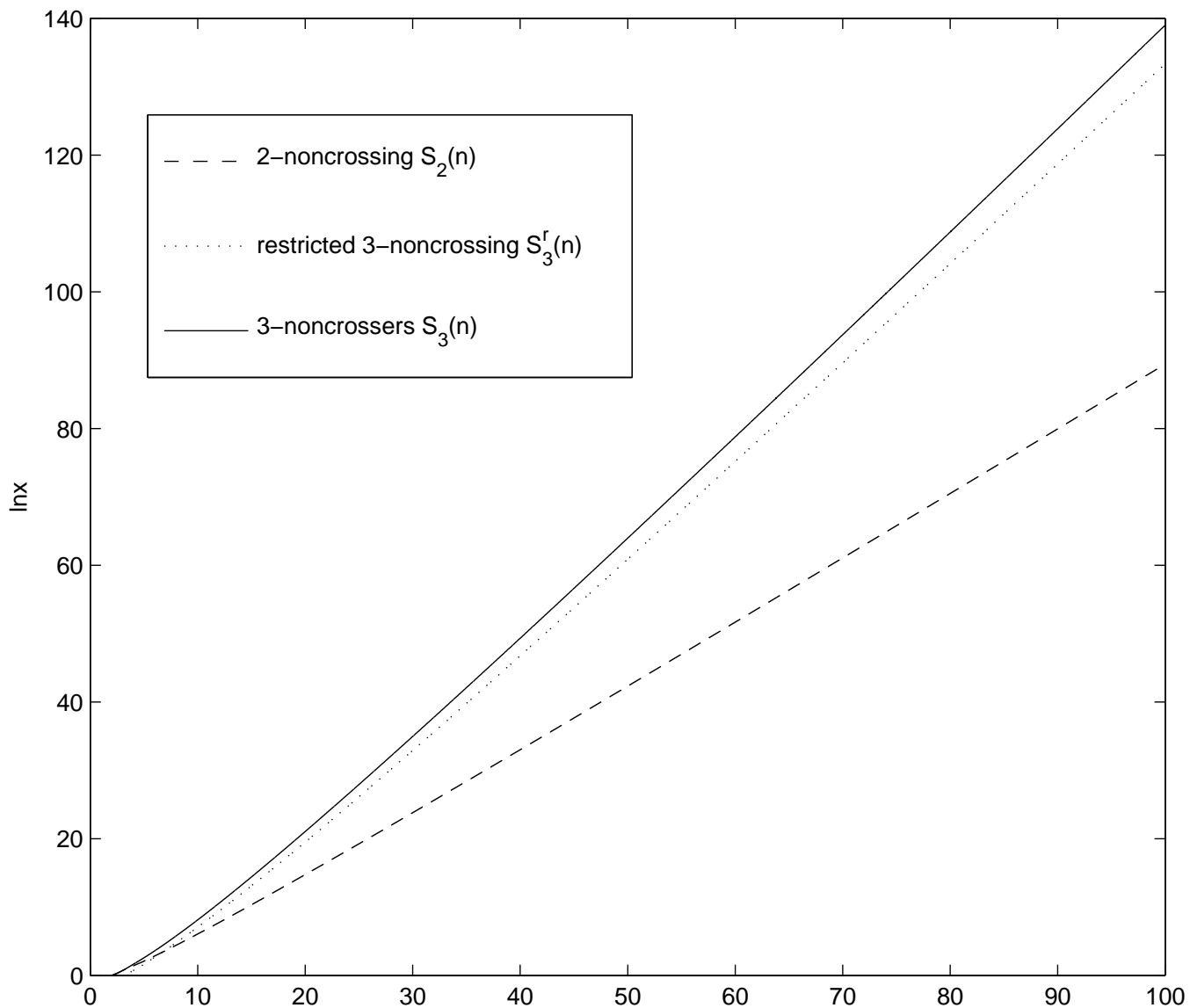# Asymptotic numbers of $3$-noncrossing RNA structures



FIGURE 10. The numbers of RNA structures for large $n$. 2-noncrossing RNA structures, 3-noncrossing RNA structures and restricted 3-noncrossing RNA structures. Numerically exponential growth rates: $S_2(n) \sim 2.5913^n$ ($n = 1000$), $S_3(n) \sim 4.6542^n$ ($n = 1000$), and $S_3^{(r)}(n) \sim 4.2741^n$ ($n = 400$).
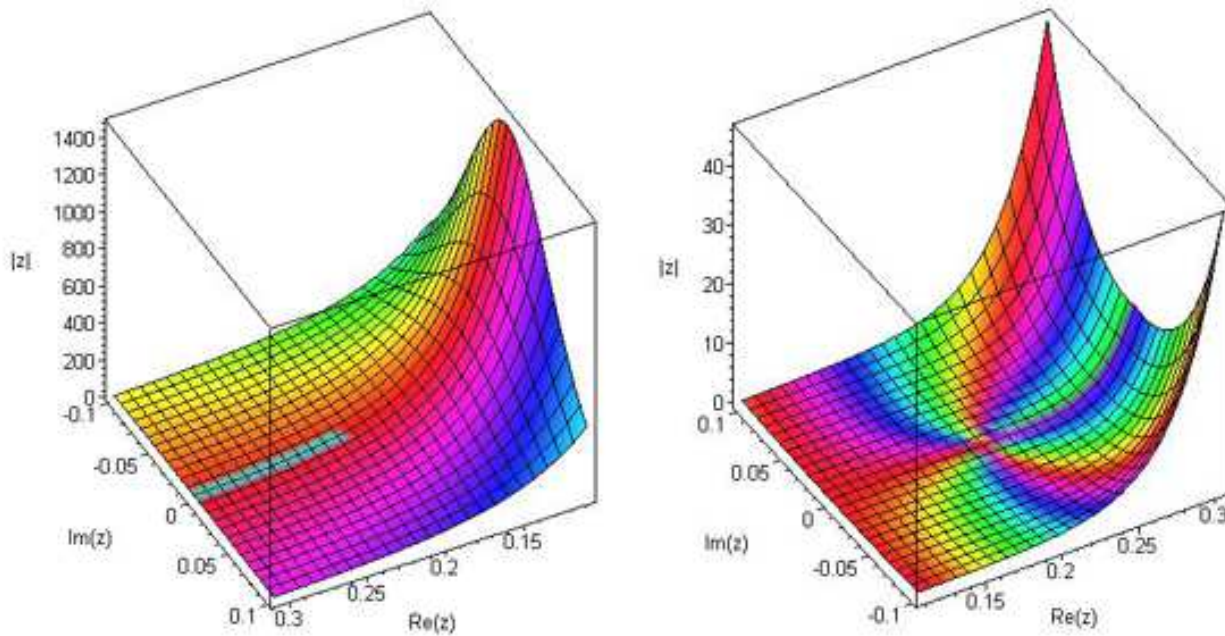
# Asymptotic Combinatorics: Toroidal Harmonics



FIGURE 11. Toroidal harmonics and its singular expansion. We display the analytic continuation of $\sum_{n\geq 0} S_3(n)z^n$, the generating function of 3-noncrossing RNA structures (left) and its singular expansion (right) at the dominant singularity $\rho_3 = \frac{5-\sqrt{21}}{2}$.

# Asymptotic Combinatorics: Toroidal Harmonics

**Lemma 1.** *Let $z$ be an indeterminant over $\mathbb{R}$ and $w \in \mathbb{R}$ a parameter. Let furthermore $\rho_k(w)$ denote the radius of convergence of the power series $\sum_{n \geq 0}[\sum_{h \leq n/2} S_k(n,h)w^{2h}]z^n$. Then for $|z| < \rho_k(w)$ holds*

$$(0.7) \qquad \sum_{n \geq 0}\sum_{h \leq n/2} S'_k(n,h)w^{2h}z^n = \frac{1}{w^2z^2 - z + 1} \sum_{n \geq 0} f_k(2n,0) \left( \frac{wz}{w^2z^2 - z + 1} \right)^{2n} .$$

*In particular we have for $w = 1$,*

$$(0.8) \qquad \sum_{n \geq 0} S_k(n)z^n = \frac{1}{z^2 - z + 1} \sum_{n \geq 0} f_k(2n,0) \left( \frac{z}{z^2 - z + 1} \right)^{2n} .$$

**Theorem 3.** *The number of $3$-noncrossing RNA structures is asymptotically given by*

$$S_3(n) \quad \sim \quad \frac{10.4724 \cdot 4!}{n(n-1)\ldots(n-4)} \left( \frac{5 + \sqrt{21}}{2} \right)^n .$$

Emma Y. Jin and Christian M. Reidys *Asymptotics of RNA Structures with Pseudoknots*, Bulletin of Math. Bio., 2007, accepted.

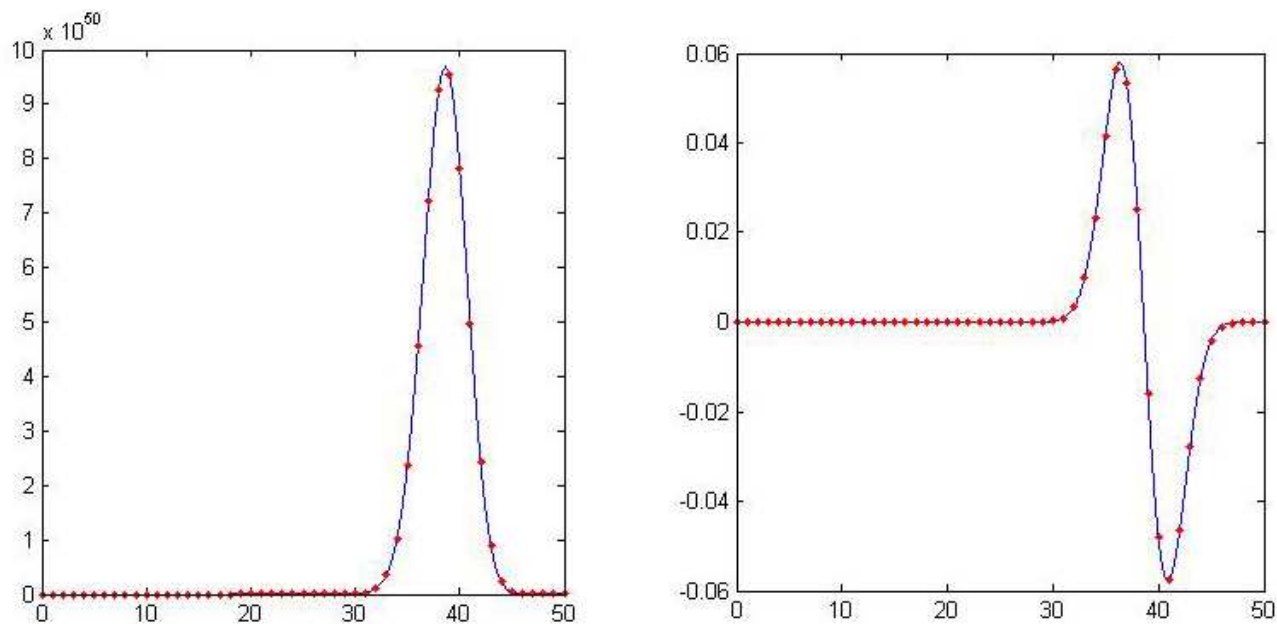# Central and Local Limit Theorems for RNA structures



FIGURE 12. Central limit theorem and local limit theorem for 3-noncrossing RNA structures of length $n = 100$ with exactly $h$ arcs: we display the central limit theorem (left) for $S'_3(100, h), h = 1, 2, \cdots 50$ (labeled by red dots) with mean $0.39089 \cdot 100 = 39.089$ and variance $0.041565 \cdot 100 = 4.1565$, and for the local limit theorem (right), we display the difference $\sqrt{4.1565} \; \mathbb{P}\left(\frac{X_n - 39.089}{\sqrt{4.1565}} = x\right) - \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$ which is maximal close to the peak of the distribution.

# Central and Local Limit Theorems for RNA structures

**Theorem 4. (Central Limit Theorem)** *Let* $S_3'(n,h)$ *be the number of* 3-*noncrossing RNA structures with exactly* $h$ *arcs. Let* $X_n$ *be the r.v. having the distribution*

(0.9) $$\forall\, h = 0, 1, \ldots \lfloor \frac{n}{2} \rfloor, \qquad \mathbb{P}(X_n = h) = \frac{S_3'(n,h)}{S_3(n)}$$

*Then the random variable* $\frac{X_n - \mu n}{\sqrt{\sigma^2 n}}$ *has asymptotically normal distribution with parameter* $(0,1)$, *i.e.*

(0.10) $$\lim_{n \to \infty} \mathbb{P}\left( \frac{X_n - \mu n}{\sqrt{\sigma^2 n}} < x \right) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{x} e^{-\frac{1}{2}t^2} dt$$

*and* $\mu, \sigma^2$ *are given by*

(0.11) $$\mu = -\frac{-\frac{3}{2} + \frac{13}{42}\sqrt{21}}{\frac{5}{2} - \frac{1}{2}\sqrt{21}} = 0.39089 \quad \text{and} \quad \sigma^2 = \mu^2 - \frac{1 - \frac{94}{441}\sqrt{21}}{\frac{5-\sqrt{21}}{2}} = 0.041565 \ .$$

**Theorem 5. (Local Limit Theorem)** *Let* $S_3'(n,h)$ *be the number of* 3-*noncrossing RNA structures with exactly* $h$ *arcs. Let* $X_n$ *be the r.v. having the distribution*

(0.12) $$\forall\, h = 0, 1, \ldots \lfloor \frac{n}{2} \rfloor, \qquad \mathbb{P}(X_n = h) = \frac{S_3'(n,h)}{S_3(n)}$$

*Then we have for set* $S = \{x \mid x = o(\sqrt{n})\}$

(0.13) $$\lim_{n \to \infty} \sup_{x \in S} \left| \sqrt{\sigma^2 n}\, \mathbb{P}\left( \frac{X_n - n\mu}{\sqrt{\sigma^2 n}} = x \right) - \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \right| = 0 \ ,$$

*where* $\mu = 0.39089$ *and* $\sigma^2 = 0.041565$.

Emma Y. Jin and Christian M. Reidys *Central and Local Limit Theorems of RNA Stuctures*, Journal of theor. Bio., 2007, submitted
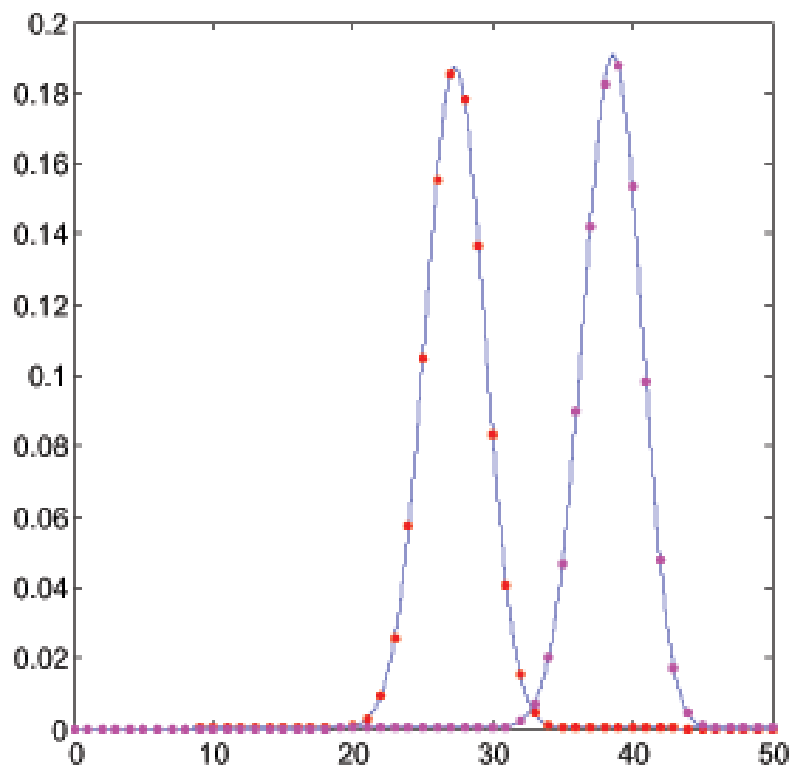
FIGURE 13. Central limit theorem of 2-noncrossing and 3-noncrossing RNA structures: both random variables are normalized to $S_2'(n,h)/S_2(n)$ and $S_3'(n,h)/S_3(n)$, respectively. In case of $n = 100$, for 2-noncrossing RNA structures we have a mean of $0.276393\,n = 27.6393$ and variance $0.044721\,n = 4.4721$ (left curve), while for 3-noncrossing RNA structures mean $0.39089\,n = 39.089$ and variance $0.041565\,n = 4.1565$ (right curve). The red dots and magenta dots represent the values $S_2'(n,h)/S_2(n)$ and $S_3'(n,h)/S_3(n)$, respectively.

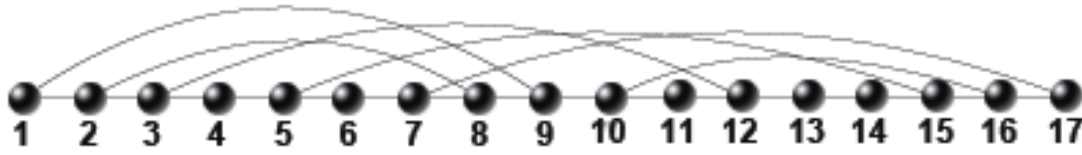# Folding of $k$-noncrossing RNA structures



FIGURE 14. A 5-noncrossing structure corresponding to the oscillating tableau below and subsequently the corresponding walk $\gamma_{a,a}$ in $\mathbb{Z}^4$.