

Distribution of motifs generated by a stochastic context-free grammar for RNA folding

Svetlana Poznanović

Department of Mathematical Sciences
Clemson University

Folding of RNA

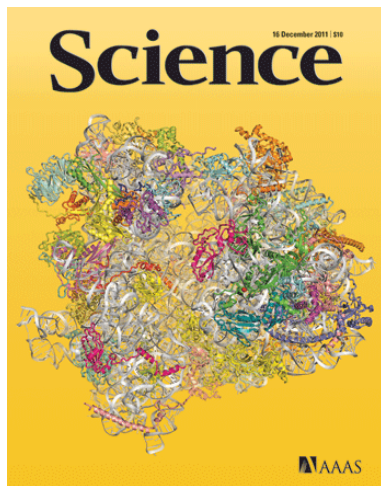
Sequence



Structure

```
GGGCGUAUGGCG  
CGUAGUCGGUAG  
CGCGCUCCCUUC  
GCCUGGGAGACU  
CCGGUGUCCGG  
ACACGUCCACCA
```



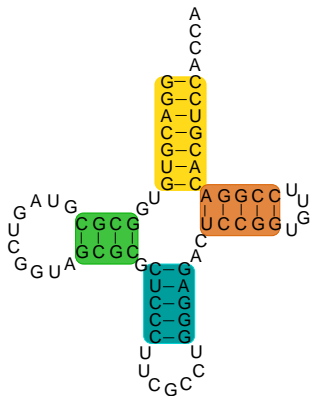


- Structure is often important for function
- Experimental determination of the structure is too often nontrivial

Intermediary structure: secondary structure

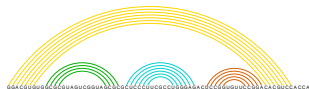
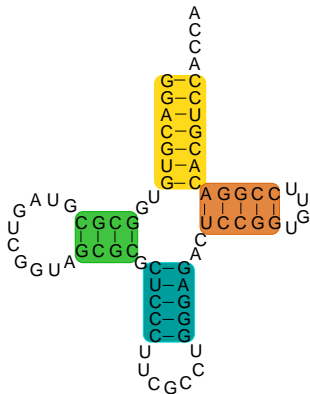
Sequence → Secondary structure → 3D structure

```
GGGCGUAUGGCG  
CGUAGUCGGUAG  
CGCGCUCCUUC  
GCCUGGGAGACU  
CCGGUGUCCGG  
ACACGUCCACCA
```



Secondary structure

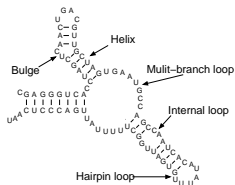
- Prediction of 3D structure is difficult, but the secondary structure can reveal a lot of information.
- Secondary structure = Non-crossing matching



Secondary structure prediction methods

Physics-based methods

- Discrete optimization problem
- Model uses thousands of measured thermodynamic parameters



Knowledge-based methods

- Parameters obtained through Maximum Likelihood Estimation

Stochastic context-free grammar (SCFG)

$$G = (V, \alpha, S, R, P_R)$$

- V is a nonterminal alphabet
- α is a terminal alphabet
- $S \in V$ is a special start symbol
- R is a set of productions in which the left hand-side is a nonterminal symbol
- P_R is a probability distribution on productions:

$$\sum_{\lambda} P(X \rightarrow \lambda) = 1, \quad \forall X \in V.$$

The Pfold SCFG

- Pfold – program for RNA secondary structure prediction which uses a SCFG (Knudsen, Hein, 1999, 2003)
- Nonterminal symbols = $\{S, L, F\}$
Terminal symbols = $\{d, d', t\}$

$$S \rightarrow LS \quad (p_1) \quad \text{or} \quad L \quad (q_1)$$

$$L \rightarrow dFd' \quad (p_2) \quad \text{or} \quad t \quad (q_2)$$

$$F \rightarrow dFd' \quad (p_3) \quad \text{or} \quad LS \quad (q_3)$$

- $p_i, q_i > 0$ – probabilities
- $p_i + q_i = 1$

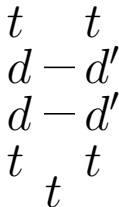
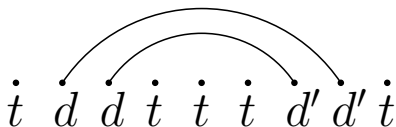
Left-to-right parse \longleftrightarrow Secondary structure

$S \xrightarrow{p_1} LS \xrightarrow{q_2} tS \xrightarrow{p_1} tLS \xrightarrow{p_2} tdFd'S \xrightarrow{p_3} tddFd'd'S \xrightarrow{q_3}$

$tddLSd'd'S \xrightarrow{q_2} tddtS d'd'S \xrightarrow{p_1} tddtLSd'd'S \xrightarrow{q_2} tddttSd'd'S \xrightarrow{q_1}$

$tddttLd'd'S \xrightarrow{q_2} tddtttd'd'S \xrightarrow{q_1} tddtttd'd'L \xrightarrow{q_2} tddtttd'd't$

- d, d' – paired nucleotides
- t – unpaired nucleotide



$S \rightarrow LS \quad (p_1) \quad \text{or} \quad L \quad (q_1)$

$L \rightarrow dFd' \quad (p_2) \quad \text{or} \quad t \quad (q_2)$

$F \rightarrow dFd' \quad (p_3) \quad \text{or} \quad LS \quad (q_3)$

“Apparently without considering alternative designs, Knudsen and Hein had already described the simple and effective G6 grammar, and extended it to analysis of input multiple sequence alignments in the program Pfold [. . .] Because our exploration has not been exhaustive, we can not say that G6 is the best possible simple grammar. However, after exploring various alternative SCFG designs, **we confirm that the Knudsen/Hein grammar is an excellent, simple framework in which to develop some probabilistic RNA analysis methods.**”

BMC Bioinformatics



Research article

Open Access

Evaluation of several lightweight stochastic context-free grammars for RNA secondary structure prediction

Robin D Dowell and Sean R Eddy*

Central limit theorems

- \mathbb{X}_n – number of occurrences of a fixed biologically interesting motif in a secondary structure with n nucleotides
- Biologically interesting motifs: base pairs, helices, multi-branch loops, hairpin loops, internal loops, bulges

Theorem (Heitsch, P.)

For a generic choice of the grammar probabilities,

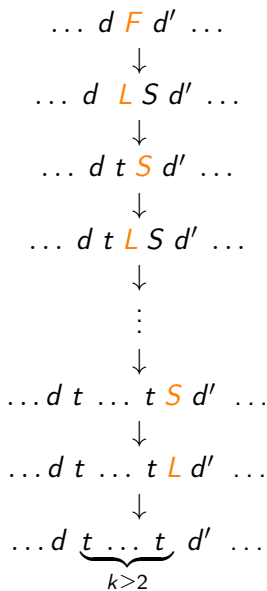
$$\mathbb{X}_n^* = \frac{\mathbb{X}_n - \mu n}{\sqrt{n\sigma^2}}$$

converge in distribution to a Gaussian variable, i.e.,

$$\lim_{n \rightarrow \infty} \mathbb{P}(\mathbb{X}_n^* < x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{c^2}{2}} dc.$$

The constants μ and σ can be explicitly computed as functions of the probabilities.

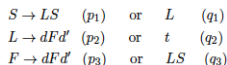
Proof sketch: Hairpin loops



- Probability for this sequence of steps:

$$q_1 q_2^2 q_3 (p_1 q_2)^{k-2}.$$

- Hairpin loops created: 1
- Nucleotides created: k



Proof sketch: From grammar to generating functions

$$\begin{cases} S \rightarrow L \text{ or } LS \\ L \rightarrow dFd' \text{ or } t \\ F \rightarrow dFd' \text{ or } LS \end{cases}$$

$$S(z, u) = \sum_{S \xrightarrow{*} M} p(M) z^{\#\text{nucleotides}(M)} u^{\#\text{hairpin loops}(M)}$$

Proof sketch: From grammar to generating functions

$$\begin{cases} S \rightarrow L \text{ or } LS \\ L \rightarrow dFd' \text{ or } t \\ F \rightarrow dFd' \text{ or } LS \end{cases}$$

$$S(z, u) = \sum_{S \xrightarrow{*} M} p(M) z^{\#\text{nucleotides}(M)} u^{\#\text{hairpin loops}(M)}$$

$$\begin{cases} S(z, u) = p_1 L(z, u) + q_1 L(z, u) S(z, u) \\ L(z, u) = p_2 z^2 F(z, u) + q_2 z \\ F(z, u) = p_3 z^2 F(z, u) + q_3 u \frac{q_1 q_2^2 z^2}{1 - p_1 q_2 z} + q_3 \left(L(z, u) S(z, u) - \frac{q_1 q_2^2 z^2}{1 - p_1 q_2 z} \right) \end{cases}$$

Proof sketch: Singularity analysis

$$S(z, u) = Q(z, u) - \frac{\sqrt{P(z, u)}}{2p_2q_3z^2(1 - p_1q_2z)}$$

$Q(z, u)$ – rational function

$P(z, u)$ – polynomial

Theorem (Flajolet, Sedgewick)

Let $G(z, u)$ be a function that is bivariate analytic at $(z, u) = (0, 0)$ and has non-negative coefficients and let \mathbb{X}_n be a random variable such that

$$\mathbb{P}(\mathbb{X}_n = k) = \frac{[z^n u^k]G(z, u)}{[z^n]G(z, 1)}.$$

Under certain technical conditions, there exist constants μ and σ such that the normalized random variable

$$\mathbb{X}_n^* = \frac{\mathbb{X}_n - \mu n}{\sqrt{n\sigma^2}}$$

converges in distribution to a Gaussian variable with a speed of convergence $O(n^{-1/2})$.

Expected number of motifs

$$\mathbb{E}(\text{Base pairs}) \sim \frac{\alpha}{\gamma} n$$

$$\mathbb{E}(\text{Helices}) \sim \frac{\alpha\beta}{\gamma} n$$

$$\mathbb{E}(\text{Hairpins}) \sim \frac{\alpha\beta(1 + p_1 q_2 \rho_0)}{4\gamma} n$$

$$\mathbb{E}(\text{Right bulges}) \sim \frac{\alpha^2\beta}{4\gamma} n$$

$$\mathbb{E}(\text{Internal loops}) \sim \frac{p_1 q_2 \rho_0 \alpha\beta}{4\gamma} n$$

$$\mathbb{E}(\text{Multi-branch loops}) \sim \frac{\alpha\beta}{4\gamma} n$$

$$\alpha = 1 - p_1 q_2 \rho_0, \beta = 1 - p_3 \rho_0^2, \gamma = 3 - p_1 q_2 \rho_0 - p_3 \rho_0^2 - p_1 p_3 q_2 \rho_0^3$$

ρ_0 is the smallest positive root of $(1 - p_1 q_2 z)^2 (1 - p_3 z^2) - 4 p_2 q_1 q_2 q_3 z^3$

Corollary

If $p_i, q_i > 0$, then

- (i) $\mathbb{E}(X_n^{lb}) = \mathbb{E}(X_n^{rb}),$
- (ii) $\mathbb{E}(X_n^m) = \frac{1}{4}\mathbb{E}(X_n^{hel})(1 + o(n)),$
- (iii) $\mathbb{E}(X_n^{hp}) = (\mathbb{E}(X_n^i) + \mathbb{E}(X_n^m))(1 + o(n)),$
- (iv) $\mathbb{E}(X_n^m) = (\mathbb{E}(X_n^{lb}) + \mathbb{E}(X_n^i))(1 + o(n)),$

Prediction using CYK

	No. Sequences (Type)	Av. Length	St. Dev. Length
Set I	122 (5S)	121.17	3.1
Set II	37 (16S)	956.46	6.51
Set III	81 (16S)	1521.33	24.86
Set IV	50 (16S)	1787.1	20.09
Set V	34 (23S)	2912.85	23.08

The transmission probabilities are

$$p_1 = 0.868534, p_2 = 0.105397, p_3 = 0.787640,$$

$$q_1 = 0.131466, q_2 = 0.894603, q_3 = 0.212360$$

and the emission probabilities are

A	0.001167	0.177977	0.001058	0.001806	A	0.364097
U	0.177977	0.002793	0.049043	0.000763	U	0.273013
G	0.001058	0.049043	0.000406	0.266974	G	0.211881
C	0.001806	0.000763	0.266974	0.000391	C	0.151009

Comparison with native secondary structures

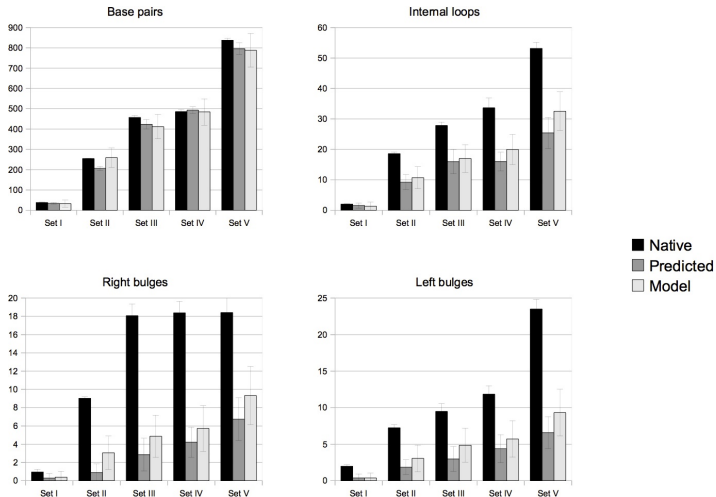


Figure : Motifs where there is agreement between the model expectations and the observed distribution in the most probable predicted structures.

Comparison with native secondary structures

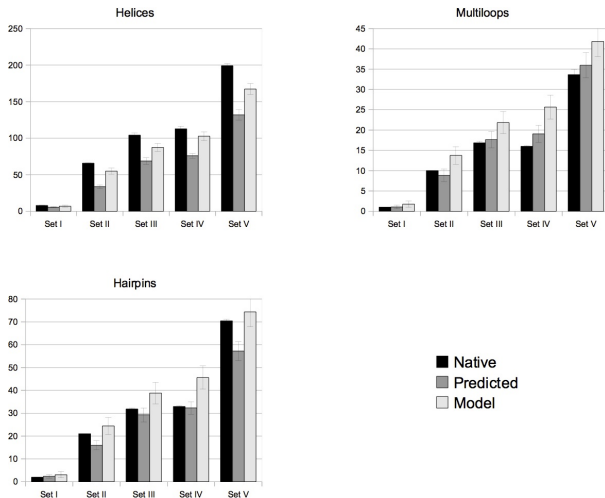


Figure : Motifs for which there is a larger discrepancy between the predicted structures and the model.

Prediction using CYK

		Ratios of Averages				
		$\frac{RB}{LB}$	$\frac{Hel}{ML}$	$\frac{IL+ML}{HP}$	$\frac{IL+LB}{ML}$	$\frac{IL+RB}{ML}$
Set I	Native	0.48	7.96	1.51	4.01	2.97
	Predicted	0.77	5.27	1.12	1.80	1.72
Set II	Native	1.25	6.58	1.36	2.58	2.76
	Predicted	0.50	3.80	1.13	1.24	1.14
Set III	Native	2.50	6.18	1.40	2.22	2.72
	Predicted	0.96	3.90	1.15	1.07	1.07
Set IV	Native	1.93	7.06	1.51	2.85	3.25
	Predicted	0.96	3.99	1.09	1.07	1.06
Set V	Native	0.78	5.92	1.23	2.28	2.13
	Predicted	1.02	3.67	1.07	0.89	0.89
Model		1	4	1	1	1

- What is the effect of changing the emission probabilities?
- (with Heitsch and Greenwood) Why does this grammar give such nice ratios?

Example:

(slightly modified KH99 - strong performer (Anderson et al.))

$$S \rightarrow LS|.|(F)$$
$$L \rightarrow .|(F)$$
$$F \rightarrow LS|(F)$$

For this grammar, the ratio multiloops/helices is ≥ 4 .

Thank you.