

# Asymptotics of Canonical and Saturated RNA Secondary Structures

Peter Clote\*    Evangelos Kranakis†    Danny Krizanc‡    Bruno Salvy§

July 9, 2009

## Abstract

It is a classical result of Stein and Waterman that the asymptotic number of RNA secondary structures is  $1.104366 \cdot n^{-3/2} \cdot 2.618034^n$ . In this paper, we study combinatorial asymptotics for two special subclasses of RNA secondary structures – *canonical* and *saturated* structures. Canonical secondary structures are defined to have no lonely (isolated) base pairs. This class of secondary structures was introduced by Bompfünnewerer et al., who noted that the run time of Vienna RNA Package is substantially reduced when restricting computations to canonical structures. Here we provide an explanation for the speed-up, by proving that the asymptotic number of canonical RNA secondary structures is  $2.1614 \cdot n^{-3/2} \cdot 1.96798^n$  and that the expected number of base pairs in a canonical secondary structure is  $0.31724 \cdot n$ . The asymptotic number of canonical secondary structures was obtained much earlier by Hofacker, Schuster and Stadler using a different method.

Saturated secondary structures have the property that no base pairs can be added without violating the definition of secondary structure (i.e. introducing a pseudoknot or base triple). Here we show that the asymptotic number of saturated structures is  $1.07427 \cdot n^{-3/2} \cdot 2.35467^n$ , the asymptotic expected number of base pairs is  $0.337361 \cdot n$ , and the asymptotic number of saturated stem-loop structures is  $0.323954 \cdot 1.69562^n$ , in contrast to the number  $2^{n-2}$  of (arbitrary) stem-loop structures as classically computed by Stein and Waterman. Finally, we apply work of Drmota [8, 9] to show that the density of states for [all resp. canonical resp. saturated] secondary structures is asymptotically Gaussian. We introduce a stochastic greedy method to sample random saturated structures, called *quasi-random saturated structures*, and show that the expected number of base pairs of is  $0.340633 \cdot n$ .

---

\*Department of Biology, Boston College, Chestnut Hill, MA 02467, USA.

Research partially supported by National Science Foundation Grants DBI-0543506, DMS-0817971, and the RNA Ontology Consortium. Additional support is gratefully acknowledged to the Foundation Digiteo-Triangle de la Physique and to Deutscher Akademischer Austauschdienst.

†School of Computer Science, Carleton University, K1S 5B6, Ottawa, Ontario, Canada. Research supported in part by Natural Sciences and Engineering Research Council of Canada (NSERC) and Mathematics of Information Technology and Complex Systems (MITACS).

‡Department of Mathematics, Wesleyan University, Middletown CT 06459, USA.

§Algorithms Project, Inria Paris-Rocquencourt, France. Supported in part by the Microsoft Research-Inria Joint Centre.

# 1 Introduction

Imagine an undirected\* graph, described by placing graph vertices  $1, \dots, n$  along the periphery of a circle in a counter-clockwise manner, and placing graph edges as chords within the circle. An *outerplanar* graph is a graph whose circular representation is planar; i.e. there are no crossings. An RNA secondary structure, formally defined in Section 2, is an outerplanar graph (no pseudoknots) with the property that no vertex is incident to more than one edge (no base triples) and that for every chord between vertices  $i, j$ , there exist at least  $\theta = 1$  many vertices that are not incident to any edge (hairpin requirement). RNA secondary structure is equivalently defined to be a well-balanced parenthesis expression  $s_1, \dots, s_n$  with dots, where if nucleotide  $i$  is unpaired then  $s_i = \bullet$ , while if there is a base pair between nucleotides  $i < j$  then  $s_i = ($  and  $s_j = )$ . This latter representation is known as the *Vienna representation* or *dot bracket notation* (dbn).

Formally, a well-balanced parenthesis expression  $w_1 \cdots w_n$  can be defined as follows. If  $\Sigma$  denotes a finite alphabet, and  $\alpha \in \Sigma$ , and  $w = w_1 \cdots w_n \in \Sigma^*$  is an arbitrary *word*, or sequence of characters drawn from  $\Sigma$ , then  $|w|_\alpha$  designates the number of occurrences of  $\alpha$  in  $w$ . Letting  $\Sigma = \{ (, ) \}$ , a word  $w = w_1 \cdots w_n \in \Sigma^*$  is *well-balanced* if for all  $1 \leq i < n$ ,  $|w_1 \cdots w_i|_{(} \geq |w_1 \cdots w_i|_{)}$  and  $|w_1 \cdots w_n|_{(} = |w_1 \cdots w_n|_{)}$ . Finally, when considering RNA secondary structures, we consider instead the alphabet  $\Sigma = \{ (, ), \bullet \}$ , but otherwise the definition of well-balanced expression remains unchanged. The number of well-balanced parenthesis expressions of length  $n$  over the alphabet  $\Sigma = \{ (, ) \}$  is known as the Catalan number  $C_n$ , while that over the alphabet  $\Sigma = \{ (, ), \bullet \}$  is known as the Motzkin number  $M_n$  [7]. Stein and Waterman [25] computed the number  $S_n$  of well-balanced parenthesis expressions in the alphabet  $\Sigma = \{ (, ), \bullet \}$ , where there exist at least  $\theta = 1$  occurrences of  $\bullet$  between corresponding left and right parentheses  $($  respectively  $)$ . It follows that  $S_n$  is exactly the number of RNA secondary structures on  $[1, n]$ , where there exist at least  $\theta = 1$  unpaired bases in every hairpin loop.

In this paper, we are interested in specific classes of secondary structure: *canonical* and *saturated* structures. A secondary structure is canonical [1] if it has no lonely (isolated) base pairs. A secondary structure is saturated [30] if no base pairs can be added without violating the notion of secondary structure, formally defined in Section 2. In order to compute parameters like asymptotic value for number of structures, expected number of base pairs, etc. throughout this paper, we adopt the model of Stein and Waterman [25]. In this model, any position (nucleotide, also known as base) can pair with any other position, and every hairpin loop must contain at least  $\theta = 1$  unpaired bases; i.e. if  $i, j$  are paired, then  $j - i > \theta$ . This latter condition is due to steric constraints for RNA. At the risk of additional effort, the combinatorial methods of this paper could be applied to handle the situation of most secondary structure software, which set  $\theta = 3$ .

## 1.1 Examples of secondary structure representations

Figure 1 gives equivalent views of the secondary structure of 5S ribosomal RNA with GenBank accession number NC\_000909 of the methane-generating archaeobacterium *Methanocaldococcus jannaschii*, as determined by comparative sequence analysis and taken from the *5S Ribosomal RNA Database* [26] located at <http://rose.man.poznan.pl/5SData/>. The sequence and its secondary structure in (Vienna) dot bracket notation are as follows:

---

\*We often describe the graph edges of an undirected graph as  $(i, j)$ , where  $i < j$ , rather than  $\{i, j\}$ .



In Section 3, we consider a natural stochastic process to generate random saturated structures, called in the sequel *quasi-random saturated structures*. The stochastic process adds base pairs, one at a time, according to the uniform distribution, without violating any of the constraints of a structure. The main result of this section is that asymptotically, the expected number of base pairs in quasi-random saturated structures is  $0.340633 \cdot n$ , rather close to the expected number  $0.337361 \cdot n$  of base pairs of saturated structures. The numerical proximity of these two values suggests that stochastic greedy methods might find application in other areas of random graph theory. In Section 4 we provide some concluding remarks. Finally, in the Appendix, we prove some technical results concerning expected stem length and the number of external loops of quasi-random saturated structures defined by a different stochastic process, distinguished by considering the uniform or Zipf probability distributions.

At the web site <http://bioinformatics.bc.edu/clotelab/SUPPLEMENTS/RNAasymptoticsCanonicalStr/> we have placed Python programs and Mathematica code used in computing and checking the asymptotic number of canonical and saturated secondary structures.

## 2 DSV methodology

In this section, we describe a combinatorial method sometimes called *DSV methodology*, after Delest, Schützenberger and Viennot, which is a special case of what is called the *symbolic method* in combinatorics, described at length in [24]. See also the Appendix of [17] for a detailed presentation of this method. This method enables one to obtain information on the number of combinatorial configurations defined by finite rules, for any size. This is done by translating those rules into equations satisfied by various *generating functions*. A second step is to extract asymptotic expansions from these equations. This is done by studying the singularities of these generating functions viewed as analytic functions.

Since our goal is to derive asymptotic numbers of structures, following standard convention we define an RNA secondary structure on a length  $n$  sequence to be a set of ordered pairs  $(i, j)$ , such that  $1 \leq i < j \leq n$  and the following are satisfied.

1. *Nonexistence of pseudoknots*: If  $(i, j)$  and  $(k, \ell)$  belong to  $S$ , then it is not the case that  $i < k < j < \ell$ .
2. *No base triples*: If  $(i, j)$  and  $(i, k)$  belong to  $S$ , then  $j = k$ ; if  $(i, j)$  and  $(k, j)$  belong to  $S$ , then  $i = k$ .
3. *Threshold requirement*: If  $(i, j)$  belongs to  $S$ , then  $j - i > \theta$ , where  $\theta$ , generally taken to be equal to 3, is the minimum number of unpaired bases in a hairpin loop; i.e. there must be at least  $\theta$  unpaired bases in a hairpin loop.

Note that the definition of secondary structure does not mention nucleotide identity – i.e. we do *not* require base-paired positions  $(i, j)$  to be occupied by Watson-Crick or wobble pairs. For this reason, at times we may say that  $S$  is a secondary structure on  $[1, n]$ , rather than saying that  $S$  is a structure for RNA sequence of length  $n$ . In particular, an expression such as “the asymptotic number of structures is  $f(n)$ ” means that the asymptotic number of structures on  $[1, n]$  is  $f(n)$ .

## Grammars

We now proceed with basic definitions related to context-free grammars. If  $A$  is a finite alphabet, then  $A^*$  denotes the set of all finite sequences (called *words*) of characters drawn from  $A$ . Let  $\Sigma$  be the set consisting of the symbols for left parenthesis  $($ , right parenthesis  $)$ , and dot  $\bullet$ , used to represent a secondary structure in Vienna notation. A *context-free* grammar (see, e.g., [15]) for RNA secondary structures is given by  $G = (V, \Sigma, \mathcal{R}, S_0)$ , where  $V$  is a finite set of nonterminal symbols (also called variables),  $\Sigma = \{\bullet, (, )\}$ ,  $S_0 \in V$  is the *start* nonterminal, and

$$\mathcal{R} \subseteq V \times (V \cup \Sigma)^*$$

is a finite set of production rules. Elements of  $\mathcal{R}$  are usually denoted by  $A \rightarrow w$ , rather than  $(A, w)$ . If rules  $A \rightarrow \alpha_1, \dots, A \rightarrow \alpha_m$  all have the same left-hand side, then this is usually abbreviated by  $A \rightarrow \alpha_1 | \dots | \alpha_m$ .

If  $x, y \in (V \cup \Sigma)^*$  and  $A \rightarrow w$  is a rule, then by replacing the occurrence of  $A$  in  $xAy$  we obtain  $xwy$ . Such a derivation in one step is denoted by  $xAy \Rightarrow_G xwy$ , while the reflexive, transitive closure of  $\Rightarrow_G$  is denoted  $\Rightarrow_G^*$ . The *language* generated by context-free grammar  $G$  is denoted by  $L(G)$ , and defined by

$$L(G) = \{w \in \Sigma^* : S_0 \Rightarrow_G^* w\}.$$

For any nonterminal  $S \in V$ , we also write  $L(S)$  to denote the language generated by rules from  $G$  when using start symbol  $S$ . A derivation of word  $w$  from start symbol  $S_0$  using grammar  $G$  is a *leftmost* derivation, if each successive rule application is applied to replace the leftmost nonterminal occurring in the intermediate expression. A context-free grammar  $G$  is *non-ambiguous*, if there is no word  $w \in L(G)$  which admits two distinct leftmost derivations. This notion is important since it is only when applied to non-ambiguous grammars that the DSV methodology leads to exact counts.

For the sake of readers unfamiliar with context-free grammars, we present some examples to illustrate the previous concepts. Consider the following grammar  $G$ , which generates the collection of well-balanced parenthesis strings, including the empty string.<sup>†</sup> Define  $G = (V, \Sigma, R, S)$ , where the set  $V$  of variables (also known as nonterminals) is  $\{S\}$ , the set  $\Sigma$  of terminals is  $\{(, )\}$ , where  $S$  is the start symbol, and where the set  $R$  of rules is given by

$$S \rightarrow \epsilon | (S) | SS$$

Here  $\epsilon$  denotes the empty string. We claim that  $G$  is an ambiguous grammar. Indeed, consider the following two leftmost derivations, where we denote the order of rule applications  $r1 := S \rightarrow \epsilon$ ,  $r2 := S \rightarrow SS$ ,  $r3 := S \rightarrow (S)$ , by placing the rule designator under the arrow. Clearly the leftmost derivation

$$S \xrightarrow{r2} SS \xrightarrow{r2} SSS \xrightarrow{r3,r1} ( ) SS \xrightarrow{r3,r1} ( ) ( ) S \xrightarrow{r3,r1} ( ) ( ) ( )$$

is distinct from the leftmost derivation

$$S \xrightarrow{r2} SS \xrightarrow{r3,r1} ( ) S \xrightarrow{r2} ( ) (S) S \xrightarrow{r3,r1} ( ) ( ) S \xrightarrow{r2} ( ) ( ) (S) \xrightarrow{r1} ( ) ( ) ( )$$

<sup>†</sup>A well-balanced parenthesis string is a word over  $\Sigma = \{(, )\}$  with as many closing parentheses as opening ones and such that when reading the word from left to right, the number of opening parentheses read is always at least as large as the number of closing parentheses. RNA secondary structures can be considered to be well-balanced parenthesis strings that also contain possible occurrences of  $\bullet$ , and for which there exist at least  $\theta$  occurrences of  $\bullet$  between corresponding left and right parentheses  $($  respectively  $)$ .

Type of nonterminal	Equation for the g.f.
$S \rightarrow T \mid U$	$S(z) = T(z) + U(z)$
$S \rightarrow TU$	$S(z) = T(z)U(z)$
$S \rightarrow t$	$S(z) = z$
$S \rightarrow \varepsilon$	$S(z) = 1$

Table 1: Translation between context-free grammars and generating functions. Here,  $G = (V, \Sigma, \mathcal{R}, S_0)$  is a given context-free grammar,  $S$ ,  $T$  and  $U$  are any nonterminal symbols in  $V$ , and  $t$  is a terminal symbol in  $\Sigma$ . The generating functions for the languages  $L(S)$ ,  $L(T)$ ,  $L(U)$  are respectively denoted by  $S(z)$ ,  $T(z)$ ,  $U(z)$ .

yet both generate the same well-balanced parenthesis string. For the same reason, the grammar with rules

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

generates precisely the collection of non-empty RNA secondary structures, yet this grammar is ambiguous, and we would obtain an overcount by applying the DSV methodology. In contrast, the grammar whose rules are

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

is easily seen to be non-ambiguous and to generate all *non-empty* RNA secondary structures.

## Generating Functions

Suppose that  $G = (V, \Sigma, \mathcal{R}, S)$  is a non-ambiguous context-free grammar which generates a collection  $L(S)$  of objects (e.g. canonical secondary structures). To this grammar is associated a generating function  $S(z) = \sum_{n=0}^{\infty} s_n z^n$ , such that the  $n$ th Taylor coefficient  $[z^n]S(z) = s_n$  represents the number of objects we wish to count. In the sequel,  $s_n$  will represent the number of canonical secondary structures for RNA sequences of length  $n$ . The DSV method uses Table 1 in order to translate the grammar rules of  $\mathcal{R}$  into a system of equations for the generating functions.

## Asymptotics

In the sequel, we often compute the asymptotic value of the Taylor coefficients of generating functions by first applying the DSV methodology, then using a simple corollary of a result of Flajolet and Odlyzko [10]. That corollary is restated here as the following theorem.

**Theorem 1 (Flajolet and Odlyzko)** *Assume that  $S(z)$  has a singularity at  $z = \rho > 0$ , is analytic in the rest of the region  $\Delta \setminus \{1\}$ , depicted in Figure 2, and that as  $z \rightarrow \rho$  in  $\Delta$ ,*

$$S(z) \sim K(1 - z/\rho)^\alpha. \tag{1}$$

*Then, as  $n \rightarrow \infty$ , if  $\alpha \notin 0, 1, 2, \dots$ ,*

$$s_n \sim \frac{K}{\Gamma(-\alpha)} \cdot n^{-\alpha-1} \rho^{-n}.$$

It is a consequence of Table 1 that the generating series of context-free grammars are algebraic (this is the celebrated theorem of Chomsky and Schützenberger [3]). In particular this implies that they have positive radius of convergence, a finite number of singularities, and their behaviour in the neighborhood of their singularities is of the type (1). (See [24, §VII.6–9] for an extensive treatment.)

A singularity of minimal modulus as in Theorem 1 is called a *dominant singularity*. The location of the dominant singularity may be a source of difficulty. The simple case is when an explicit expression is obtained for the generating functions; this happens for canonical secondary structures. The situation when only the system of polynomial equations is available is more involved; we show how to deal with it in the case of saturated structures.

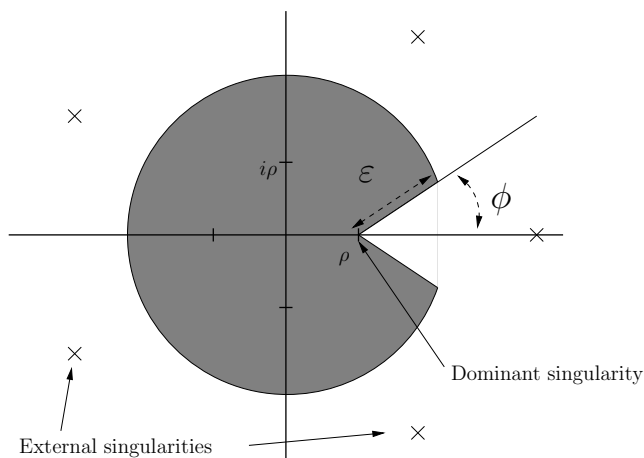


Figure 2: The shaded region  $\Delta$  where, except at  $z = \rho$ , the generating function  $S(z)$  must be analytic.

## 2.1 Asymptotic number of canonical secondary structures

In Bompfünowerer et al. [1], the notion of *canonical secondary structure*  $S$  is defined as a secondary structure having no *lonely* (isolated) base pairs; i.e. formally, there are no base pairs  $(i, j) \in S$  for which both  $(i - 1, j + 1) \notin S$  and  $(i + 1, j - 1) \notin S$ . In this section, we compute the asymptotic number of canonical secondary structures. Throughout this section, secondary structure is interpreted to mean a secondary structure on an RNA sequence of length  $n$ , for which each base can pair with any other base (not simply Watson-Crick and wobble pairs), and with minimum number  $\theta$  of unpaired bases in every hairpin loop set to be 1. At the cost of working with more complex expressions, by the same method, one could analyze the case when  $\theta = 3$ , which is assumed for the software `mfold` [31] and `RNAfold` [11].

### Grammar

Consider the context-free grammar  $G = (V, \Sigma, \mathcal{R}, S)$ , where  $V$  consists of nonterminals  $S, R$ ,  $\Sigma$  consists of the terminals  $\bullet, (, )$ ,  $S$  is the start symbol and  $\mathcal{R}$  consists of the following rules:

$$\begin{aligned} S &\rightarrow \bullet | S \bullet | (R) | S(R) \\ R &\rightarrow ( \bullet ) | (R) | (S(R)) | (S \bullet ) \end{aligned} \tag{2}$$

The nonterminal  $S$  is intended to generate all *nonempty canonical* secondary structures. In contrast, the nonterminal  $R$  is intended to generate all secondary structures which become canonical when surrounded by a closing set of parentheses. We prove by induction on expression length that the grammar  $G$  is non-ambiguous and generates all nonempty canonical secondary structures.

Define context-free grammar  $G_R$  to consist of the collection  $\mathcal{R}$  of rules from  $G$ , defined above, with starting nonterminal  $S$ , respectively. Formally,

$$G_R = (V, \Sigma, \mathcal{R}, R).$$

Let  $L(G)$ ,  $L(G_R)$  denote the languages generated respectively by grammars  $G, G_R$ . Now define languages  $L_1, L_2$  of *nonempty* secondary structures with  $\theta = 1$  by

$$\begin{aligned} L_1 &= \{S : S \text{ is canonical}\} \\ L_2 &= \{S : (S) \text{ is canonical}\}. \end{aligned}$$

Note that structures like  $\bullet\bullet(\bullet)$  and  $(\bullet)(\bullet)$  belong to  $L_1$ , but not to  $L_2$ , while structures like  $((\bullet))$  belong to both  $L_1, L_2$ . Note that any structure  $S$  belonging to  $L_2$  must be of the form  $(S_0)$ ; indeed, if  $S$  were not of this form, but rather of the form either  $\bullet S_0$  or  $(S_0)S_1$ , then by  $(S)$  would have an outermost lonely pair of parentheses.

CLAIM.  $L_1 = L(G)$ ,  $L_2 = L(G_R)$ .

PROOF OF CLAIM. Clearly  $L_1 \supseteq L(G)$ ,  $L_2 \supseteq L(G_R)$ , so we show the reverse inclusions by induction; i.e. by induction on  $n$ , we prove that  $L_1 \cap \Sigma^n \subseteq L(G) \cap \Sigma^n$ ,  $L_2 \cap \Sigma^n \subseteq L(G_R) \cap \Sigma^n$ .

BASE CASE:  $n = 1$ . Clearly  $L(G) \cap \Sigma = \{\bullet\} = L_1 \cap \Sigma$ ,  $L(G_R) \cap \Sigma = \emptyset = L_2 \cap \Sigma$ .

INDUCTION CASE: Assume that the claim holds for all  $n < k$ .

*Subcase 1.* Let  $\mathcal{S}$  be a canonical secondary structure with length  $|\mathcal{S}| = k > 1$ . Then either (1)  $\mathcal{S} = \bullet S_0$ , where  $S_0 \in L_1$ , or (2)  $\mathcal{S} = (S_0)$ , where  $S_0 \in L_2$ , or (3)  $\mathcal{S} = (S_0)S_1$ , where  $S_0 \in L_2$  and  $S_1 \in L_1$ . Each of these cases corresponds to a different rule having left side  $S$ , hence by the induction hypothesis, it follows that  $\mathcal{S} \in L(G)$ .

*Subcase 2.* Let  $\mathcal{S} \in L_2$  be a secondary structure with length  $|\mathcal{S}| = k > 1$ , for which  $(\mathcal{S})$  is canonical. If  $\mathcal{S}$  were of the form  $\bullet S_0$  or  $(S_0)S_1$ , then  $(\mathcal{S})$  would not be canonical, since its outermost parenthesis pair would be a lonely pair. Thus  $\mathcal{S}$  is of the form  $(S_0)$ , where either (1)  $S_0$  begins with  $\bullet$ , or (2)  $S_0$  is of the form  $(S_1)$ , where  $S_1$  is not canonical, but  $(S_1)$  becomes canonical, or (3)  $S_0$  is of the form  $(S_1)$ , where  $S_1$  is canonical and  $(S_1)$  is canonical as well.

In case (1),  $S_0$  is either  $\bullet$  or  $\bullet S_1$ , where  $S_1$  is canonical. In case (2),  $S_0$  is of the form  $(S_1)$ , where  $S_1$  must have the property that  $(S_1)$  is canonical. In case (3),  $S_0$  is of the form  $(S_1)S_2$ , where it must be that  $(S_1)$  is canonical and  $S_2$  is canonical. By applying corresponding rules and the induction hypothesis, it follows that  $\mathcal{S} \in L(G_R)$ .

It now follows by induction that  $L_1 = L(G)$ ,  $L_2 = L(G_R)$ . A similar proof by induction shows that the grammar  $G$  is non-ambiguous.



## Generating Functions

Now, let  $s_n$  denote the number of canonical secondary structures on a length  $n$  RNA sequence. Then  $s_n$  is the  $n$ th Taylor coefficient of the generating function  $S(z) = \sum_{n \geq 0} s_n z^n$ , denoted by  $s_n = [z^n]S(z)$ . Similarly, let  $R(z) = \sum_{n \geq 0} R_n z^n$  be the generating function for the number of secondary structures on  $[1, n]$  with  $\theta = 1$ , which become canonical when surrounded by a closing set of parentheses.

By Table 1, the non-ambiguous grammar (2) gives the following equations

$$S(z) = z + S(z)z + R(z)z^2 + S(z)R(z)z^2 \quad (3)$$

$$R(z) = z^3 + R(z)z^2 + S(z)R(z)z^4 + S(z)z^3 \quad (4)$$

which can be solved explicitly (solve the second equation for  $R$  and inject this in the first equation):

$$S(z) = \frac{1 - z - z^2 + z^3 - z^5 - \sqrt{F(z)}}{2z^4} \quad (5)$$

and

$$S(z) = \frac{1 - z - z^2 + z^3 - z^5 + \sqrt{F(z)}}{2z^4} \quad (6)$$

where

$$F(z) = 4z^5 (-1 + z^2 - z^4) + (-1 + z + z^2 - z^3 + z^5)^2. \quad (7)$$

When evaluated at  $z = 0$ , Equation (6) gives  $\lim_{r \rightarrow 0} S(z) = \infty$ . Since  $S(z)$  is known to be analytic at 0, we conclude that  $S(z)$  is given by (5).

## Location of the dominant singularity

The square root function  $\sqrt{z}$  has a singularity at  $z = 0$ , so we are led to investigate the roots of  $F(z)$ . A numerical computation with Mathematica™ gives the 10 roots 0.508136, 4.11674,  $-0.868214 - 0.619448i$ ,  $-0.868214 + 0.619448i$ ,  $-0.799805 - 0.367046i$ ,  $-0.799805 + 0.367046i$ ,  $0.410134 - 0.564104i$ ,  $0.410134 + 0.564104i$ ,  $0.945448 - 0.470929i$ ,  $0.945448 + 0.470929i$ . It follows that  $\rho = 0.508136$  is the root of  $F(z)$  having smallest (complex) modulus.

## Asymptotics

Let  $T(z) = \frac{1 - z - z^2 + z^3 - z^5}{2z^4}$  and factor  $1 - z/\rho$  out of  $F(z)$  to obtain  $Q(z)(1 - z/\rho) = F(z)$ . It follows that

$$S(z) - T(\rho) = \frac{\sqrt{Q(\rho)}}{2\rho^4} \cdot (1 - z/\rho)^\alpha + O(1 - z/\rho), \quad z \rightarrow \rho,$$

where  $\alpha = 1/2$ . This shows that  $\rho$  is indeed a dominant singularity for  $S$ . Note that for each  $n \geq 1$ ,  $S(z)$  and  $S(z) - T(\rho)$  have the same Taylor coefficient of index  $n$ , namely  $s_n$ . Now, it is a direct consequence of Theorem 1 that

$$s_n \sim \frac{K(\rho)}{\Gamma(-\alpha)} \cdot n^{-\alpha-1} \cdot (1/\rho)^n, \quad n \rightarrow \infty \quad (8)$$

where  $\alpha = 1/2$  and  $K(z) = \frac{\sqrt{Q(z)}}{2z^4}$ . Plugging  $\rho = 0.508136$  into equation (8), we derive the following theorem, first obtained by Hofacker, Schuster and Stadler [12] by a different method.

**Theorem 2** *The asymptotic number of canonical secondary structures on  $[1, n]$  is*

$$2.1614 \cdot n^{-3/2} \cdot 1.96798^n. \quad (9)$$

## 2.2 Asymptotic expected number of base pairs in canonical structures

In this section, we derive the expected number of base pairs in canonical secondary structures on  $[1, n]$ .

### Generating Functions

The DSV methodology is actually able to produce *multivariate* generating series. Modifying the equations (3,4) by adding a new variable  $u$ , intended to count the number of base pairs, we get

$$S(z, u) = z + S(z, u)z + R(z, u)uz^2 + S(z, u)R(z, u)uz^2 \quad (10)$$

$$R(z, u) = uz^3 + R(z, u)uz^2 + S(z, u)R(z, u)u^2z^4 + S(z, u)uz^3. \quad (11)$$

This can be solved as before to yield the solution<sup>‡</sup>

$$\begin{aligned} S(z, u) &= \sum_{n \geq 0} \sum_{k \geq 0} s_{n,k} z^n u^k \\ &= 2u^2 z^4 \left( 1 - z - uz^2 + uz^3 - u^2 z^5 - \right. \\ &\quad \left. \sqrt{4u^2 z^5 (-1 + uz^2 - u^2 z^4) + (-1 + z + uz^2 - uz^3 + u^2 z^5)^2} \right) \end{aligned}$$

Here, the coefficient  $s_{n,k}$  is the number of canonical secondary structures of size  $n$  with  $k$  base pairs. Using a classical observation on multivariate generating functions, we recover the expected number of base pairs in a canonical secondary structure on  $[1, n]$  using the partial derivative of  $S(z, u)$ ; indeed,

$$\begin{aligned} \frac{[z^n] \frac{\partial S(z, u)}{\partial u} (z, 1)}{[z^n] S(z, 1)} &= \frac{[z^n] \left( \sum_{i \geq 0} \sum_{k \geq 0} s_{i,k} z^i k u^{k-1} \right) (z, 1)}{s_n} \\ &= \frac{\sum_{k \geq 0} s_{n,k} k}{s_n} = \sum_{k \geq 0} k \frac{s_{n,k}}{s_n}, \end{aligned}$$

and  $s_{n,k}/s_n$  is the (uniform) probability that a canonical secondary structure on  $[1, n]$  has exactly  $k$  base pairs.

We compute that  $G(z) = \frac{\partial S(z, u)}{\partial u} (z, 1)$  satisfies

$$G(z) = \frac{-(z^2 - 2)(T(z) - \sqrt{F(z)} + z\sqrt{F(z)})}{2z^4 \sqrt{F(z)}}$$

where  $T(z) = (1 - 2z + 2z^3 - z^4 - 3z^5 + z^6)$  and  $F(z)$  is as in (7). Simplification yields

$$G(z) = \frac{-(z^2 - 2)(z - 1)}{2z^4} - \frac{T(z)(z^2 - 2)}{2z^4} \cdot \left( \frac{1}{\sqrt{F(z)}} \right).$$

<sup>‡</sup>Since  $S(z, u)$  is known to be analytic at 0, we have discarded one of the two solutions as before.

## Asymptotics

From this expression, it is clear that the dominant singularity is again located at the same  $\rho = 0.508136$ . A local expansion there gives

$$G(z) \sim K(\rho)(1 - z/\rho)^{-1/2}, \quad z \rightarrow \rho$$

with  $K(z) = -\frac{Q(z)^{-1/2}T(z)(z^2-2)}{2z^4}$ . By Theorem 1, we obtain the asymptotic value

$$\frac{K(\rho)}{\Gamma(-\alpha)} \cdot n^{-3/2} \cdot (1/\rho)^n. \quad (12)$$

Plugging  $\rho = 0.508136$  into equation (12), we find the asymptotic value of  $[z^n] \frac{\partial S(z,u)}{\partial u}(z, 1)$  is

$$0.68568 \cdot n^{-1/2} \cdot 1.96798^n. \quad (13)$$

Dividing (13) by the asymptotic number  $[z^n]S(z)$  of canonical secondary structures, given in (9), we have the following theorem.

**Theorem 3** *The asymptotic expected number of base pairs in canonical secondary structures is  $0.31724 \cdot n$ .*

## 2.3 Asymptotic number of saturated structures

An RNA secondary structure is *saturated* if it is not possible to add any base pair without violating the definition of secondary structures. If one models the folding of an RNA secondary structure as a random walk on a Markov chain (i.e. by the Metropolis-Hastings algorithm), then saturated structures correspond to *kinetic traps* with respect to the Nussinov energy model [21]. The asymptotic number of saturated structures was determined in [4] by using a method known as Bender's Theorem, as rectified by Meir and Moon [19]. In this section, we apply the DSV methodology to obtain the same asymptotic limit, and in the next section we obtain the expected number of base pairs of saturated structures.

### Grammar

Consider the context-free grammar with nonterminal symbols  $S, R$ , terminal symbols  $\bullet, (, )$ , start symbol  $S$  and production rules

$$S \rightarrow \bullet | \bullet \bullet | R \bullet | R \bullet \bullet | (S) | S(S) \quad (14)$$

$$R \rightarrow (S) | R(S) \quad (15)$$

It can be shown by induction on expression length that  $L(S)$  is the set of saturated structures, and  $L(R)$  is the set of saturated structures with no *visible* position; i.e. external to every base pair [4]. Here, position  $i$  is visible in a secondary structure  $T$  if it is external to every base pair of  $T$ ; i.e. for all  $(x, y) \in T$ ,  $i < x$  or  $i > y$ .

## Generating Functions

Let

$$S(z) = \sum_{i=0}^{\infty} s_i \cdot z^i, \quad R(z) = \sum_{i=0}^{\infty} r_i \cdot z^i \quad (16)$$

denote the generating functions  $S$  resp.  $R$ , corresponding to the problems of counting number of saturated secondary structures resp. number of saturated structures having no visible positions. Applying Table 1, we are led to the equations

$$S = z + z^2 + zR + z^2R + z^2S + z^2S^2 \quad (17)$$

$$R = z^2S + z^2RS. \quad (18)$$

## Location of the dominant singularity

By first solving (18) for  $R$  and injecting in (17), we get

$$S = z + z^2 + z^2S + z^2S^2 + (z + z^2) \frac{z^2S}{1 - z^2S}, \quad (19)$$

which upon normalizing gives a polynomial equation of the third degree

$$P(z, S) = -S^3z^4 + z(1 + z) - S^2z^2(-2 + z^2) + S(-1 + z^2) = 0. \quad (20)$$

Unlike earlier work in this paper, direct solution of this equation by Cardano's formulas gives expressions that are difficult to handle. Instead, we locate the singularity by appealing to general techniques for implicit generating functions [24, §VII.4].

By the implicit function theorem, singularities of  $P(z, S)$  only occur when both  $P$  and its partial derivative

$$\frac{\partial P}{\partial S}(z, S) = -1 + (1 + 4S)z^2 - S(2 + 3S)z^4 \quad (21)$$

vanish simultaneously.

The common roots of  $P$  and  $\partial P/\partial S$  can be located by eliminating  $S$  between those two equations, for instance using the classical theory of *resultants* (see, e.g., [14]). This gives a polynomial

$$Q(z) = z^{11}(1 + z)(4 + z - 7z^2 - 28z^3 - 32z^4 + 4z^6), \quad (22)$$

that vanishes at all  $z$  such that  $(z, S)$  is a common root of  $P$  and  $\partial P/\partial S$ .

Numerical computation of the roots of  $Q$  yields  $0, -1, -2.29493, -0.854537, -0.244657 - 0.5601i, -0.244657 + 0.5601i, 0.424687, 3.2141$ .

A subtle difficulty now lies in selecting among those points the dominant singularity of the analytic continuation of the solution  $S$  of (19) corresponding to the combinatorial problem. Indeed, it is possible that one solution of (19) is singular at a given  $r$  without the solution of interest being singular there. Considering such a singularity would result in an asymptotic expansion that is wrong by an exponential factor. One way to select the correct singularity is to apply a result by Meir and Moon [18] to Equation (19). This results in a variant of the computation in [4].

Instead, we use Pringsheim's theorem (see, e.g., [24]).

**Theorem 4 (Pringsheim)** *If  $S(z)$  has a series expansion at 0 that has nonnegative coefficients and a radius of convergence  $R$ , then the point  $z = R$  is a singularity of  $S(z)$ .*

In our example, there are only two possible real positive singularities, 0.424687 and 3.2141. The latter cannot be dominant, since it would lead to asymptotics of the form  $3.2141^{-n}$ , i.e., an exponentially decreasing number of structures. Thus the dominant singularity is at  $\rho = 0.424687$ . Since the moduli of the non-real roots of  $Q$  is  $0.611203 > \rho$ , the conditions of Theorem 1 hold, provided the function behaves as required as  $z \rightarrow \rho$ .

### Asymptotics

We now compute the local expansion of  $S(z)$  at  $\rho$ . From equation (21), we have that

$$P(\rho, S) = 0.605047 - 0.819641S + 0.328189S^2 - 0.0325295S^3 \quad (23)$$

whose (numerical approximations of) roots are the double root  $S = 1.6569$  and single root  $S = 6.77518$ . It is easily checked that 1.6569 is the only root of equation (23) in which  $P(\rho, S)$  is increasing; thus we let  $T = 1.6569$ .

Recall Taylor's theorem in two variables

$$f(x, y) = \sum_{n=0}^{\infty} \sum_{k=0}^{\infty} \frac{\partial^{n+k} f(x_0, y_0)}{\partial x^n \partial y^k} \cdot \frac{(x - x_0)^n}{n!} \cdot \frac{(y - y_0)^k}{k!}.$$

We now expand  $P(z, S)$  at  $z = \rho$  and  $S = T$  and invert this expansion. This yields

$$P(z, S) = P(\rho, T) + \frac{\partial P}{\partial S}(\rho, T)(S - T) + \frac{\partial P}{\partial z}(\rho, T)(z - \rho) + \frac{1}{2} \frac{\partial^2 P}{\partial S^2}(\rho, T)(S - T)^2 + \dots \quad (24)$$

where the dots indicate terms of higher order. The first two terms are 0, so by denoting  $P_z = \frac{\partial P}{\partial z}(\rho, T)$  and  $P_{SS} = \frac{\partial^2 P}{\partial S^2}(\rho, T)$ , we have

$$0 = P = P_z(z - \rho) + \frac{1}{2} P_{zz}(S - T)^2 + O(S - T)^3 + O((z - \rho)(S - T)^2) + O((z - \rho)^2). \quad (25)$$

Isolating  $(S - T)^2$  we get

$$\begin{aligned} (S - T)^2 &= \frac{-2P_z(z - \rho)}{P_{SS}} + O((z - \rho)^2) + O((S - T)^3) \\ S - T &= \pm \sqrt{\frac{2\rho P_z}{P_{SS}}} \cdot \sqrt{1 - z/\rho} + O(z - \rho). \end{aligned}$$

Since  $[z^n]S(z)$  is the number of saturated secondary structures on  $[1, n]$  and the Taylor coefficients in the expansion of  $\sqrt{1 - z/\rho}$  are negative, we discard the positive root and thus obtain

$$S - T = -\sqrt{\frac{2\rho P_z}{P_{SS}}} \cdot \sqrt{1 - z/\rho} + O(z - \rho). \quad (26)$$

We now make use of Theorem 1 as before and recover the following result, proved earlier in [4] by the Bender-Meir-Moon method.

**Theorem 5** *The asymptotic number of saturated structures is  $1.07427 \cdot n^{-3/2} \cdot 2.35468^n$ .*

## 2.4 Expected number of base pairs of saturated structures

In this section, we compute the expected number of base pairs of saturated structures, proceeding as in Section 2.2 by first modifying the equations to obtain bivariate generating functions and then differentiating with respect to the new variable and evaluating at 1 to obtain the asymptotic expectation.

### Generating Functions

We first modify equations (17,18) by introducing the auxiliary variable  $u$ , responsible for counting the number of base pairs:

$$S = z + z^2 + zR + z^2R + uz^2S + uz^2S^2 \quad (27)$$

$$R = uz^2S + uz^2RS. \quad (28)$$

Solving the second equation for  $R$  and injecting into the first one gives

$$P(z, u, S) = Suz^2(z + z^2) - (-1 + Suz^2)(-S + z + z^2 + Suz^2 + S^2uz^2). \quad (29)$$

### Asymptotics

We are interested in the coefficients of  $\partial S/\partial u$  at  $u = 1$ . Differentiating (29) with respect to  $u$  gives

$$\frac{\partial P}{\partial u} + \frac{\partial P}{\partial S} \frac{\partial S}{\partial u} = 0.$$

Using equation (26), we replace  $S(z, 1)$  by  $T + K\sqrt{1 - z/\rho} + O(1 - z/\rho)$  in this equation to obtain

$$\left( \rho^2 T(1 + 2(1 - \rho^2)T - 2\rho^2 T^2) + O(\sqrt{1 - z/\rho}) \right) + \left( (4K\rho^2 - 2K\rho^4 - 6K\rho^4 T)\sqrt{1 - z/\rho} + O(1 - z/\rho) \right) \frac{\partial S}{\partial u} \Big|_{u=1} = 0$$

and finally

$$\frac{\partial S}{\partial u}(z, 1) \sim -\frac{0.642305}{\sqrt{1 - z/\rho}}.$$

Applying Theorem 1 to equation (30) gives

$$\rho^n [z^n] \frac{\partial S}{\partial u}(z, 1) \sim \frac{0.642305}{\Gamma(1/2)} \cdot n^{-1/2} = 0.362417 \cdot n^{-1/2}.$$

It follows that the asymptotic expected number of base pairs in saturated structures on  $[1, n]$  is

$$\frac{[z^n] \frac{\partial S(z, u)}{\partial u}(z, 1)}{[z^n] S(z, 1)} \sim \frac{0.362417 \cdot n^{-1/2} \cdot \rho^{-n}}{1.07427 \cdot n^{-3/2} \cdot \rho^{-n}} = 0.337361 \cdot n$$

We have just proved the following.

**Theorem 6** *The asymptotic expected number of base pairs for saturated structures is  $0.337361 \cdot n$ .*

Since the Taylor coefficient  $s_{n,k}$  of generating function  $S(z, u) = \sum_{n,k} s_{n,k} z^n u^k$  is equal to the number of saturated structures having  $k$  base pairs, it is possible that the methods of this section will suffice to solve the following open problem.

**Open Problem 1** *Clearly, the maximum number of base pairs in a saturated structure on  $[1, n]$  where  $\theta = 1$  is  $\lfloor \frac{n-1}{2} \rfloor$ . For fixed values of  $k$ , what is the asymptotic number  $s_{n, \lfloor (n-1)/2 \rfloor - k}$  of saturated secondary structures having exactly  $k$  base pairs fewer than the maximum?*

Note that in [4], we solved this problem for  $k = 0, 1$ .

A related interesting question concerns whether the number of secondary structures  $s_{n,k}$  having  $k$  base pairs is approximately Gaussian. As first suggested by Y. Ponty (personal communication), this is indeed the case. More formally, consider for fixed  $n$  the the finite distribution  $\mathbb{P}_n = p_1, \dots, p_n$ , where  $p_k = s_{n,k}/s_n$  and  $s_n = \sum_k s_{n,k}$ . In the Nussinov energy model, the energy of a secondary structure with  $k$  base pairs is  $-k$ , so the distribution  $\mathbb{P}_n$  is what is usually called the *density of states* in physical chemistry. It follows from Theorem 1 of of Drmota [9] (see also [8]) that  $\mathbb{P}_n$  is Gaussian. Similarly, it follows from Theorem 1 of Drmota that the asymptotic distribution of density of states of both canonical and saturated structures is Gaussian. Details of a Maple session applying Drmota's theorem to saturated structures appears in the web supplement <http://bioinformatics.bc.edu/clotelab/SUPPLEMENTS/RNAasymptoticsCanonicalStr/>.

## 2.5 Asymptotic number of saturated stem-loops

Define a *stem-loop* to be a secondary structure  $S$  having a unique base pair  $(i_0, j_0) \in S$ , for which all other base pairs  $(i, j) \in S$  satisfy the relation  $i < i_0 < j_0 < j$ . In this case,  $(i_0, j_0)$  defines a hairpin, and the remaining base pairs, as well as possible internal loops and bulges, constitute the stem. We have the following simple result due to Stein and Waterman [25].

**Proposition 1** *There are  $2^{n-2} - 1$  stem-loop structures<sup>§</sup> on  $[1, n]$ .*

**Proof.** Let  $L(n)$  denote the number of secondary structures with *at most* one loop on  $(1, \dots, n)$ . Then  $L(1) = 1 = L(2)$ . There are two cases to consider for  $L(n+1)$ .

CASE 1. If  $n+1$  does not form a base pair, then we have a contribution of  $L(n)$ .

CASE 2.  $n+1$  forms a base pair with some  $1 \leq j \leq n-1$ . In this case, since only one hairpin loop is allowed, there is no base-pairing for the subsequence  $s_1, \dots, s_{j-1}$ , and hence if  $n+1$  base-pairs with  $j$ , then we have a contribution of  $L(n - (j+1) + 1) = L(n-j)$ . Hence

$$\begin{aligned} L(n+1) &= L(n) + \sum_{j=1}^{n-1} L(n-j) \\ &= L(n) + L(n-1) + \dots + L(1) \end{aligned}$$

and hence  $L(1) = 1$ ,  $L(2) = 1$ ,  $L(3) = 2$ , and from there  $L(n) = 2^{n-2}$  by induction. ■

---

<sup>§</sup>In [25], stem-loop structures are called *hairpins*. Since the appearance of [25], common convention is that a hairpin is a structure consisting of a single base pair enclosing a loop region; i.e.  $(\bullet \cdots \bullet)$ . Here we use the more proper term *stem-loop*.

We now compute the asymptotic number of *saturated* stem-loop structures. Let  $h(n)$  be the number of saturated stem-loops on  $[1, n]$ , defined by  $h(n) = 1$  for  $n = 0, 1, 2, 3$ ,  $h(4) = 3$ , and

$$h(n) = h(n-2) + 2h(n-3) + 2h(n-4) \quad (30)$$

for  $n \geq 5$ . Note that we have defined  $h(1) = 1 = h(2)$  for notational ease in the sequel, although there are in fact no stem-loops of size 1 or 2. Indeed in this case, the only structures of size 1 respectively 2 are  $\bullet$  and  $\bullet\bullet$ .

The first few terms in the sequence  $h(1), h(2), h(3), \dots$  are 1, 1, 1, 3, 5, 7, 13, 23, 37, 63, 109, 183, 309, 527, 893, 1511, 2565, 4351, 7373, 12503; for instance,  $h(20) = 12503$ .

### Grammar

It is easily seen that the following rules

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

provide for a non-ambiguous context-free grammar to generate all non-empty saturated stem-loops. It defines actually a special kind of context-free language, called regular, whose generating function is rational.

### Generating Function

By the DSV methodology, we obtain the functional relation

$$R(z) = z + z^2 + R(z)z^2 + 2R(z)z^3 + 2R(z)z^4$$

whose solution is the rational function

$$R(z) = \frac{P(z)}{Q(z)} = \frac{z}{1 - z - 2z^3} \quad (31)$$

where  $P(z) = z$  and  $Q(z) = 1 - z - 2z^3$ .

### Asymptotics

For rational functions, an easy way to compute the asymptotic behaviour of the Taylor coefficients is to compute a partial fraction decomposition and isolate the dominant part. This is equivalent to solving the corresponding linear recurrence. See also [23, p. 325] or [22, Thm. 9.2].

Partial fraction decomposition yields

$$R(z) = \frac{A(a_1)}{1 - z/a_1} + \frac{A(a_2)}{1 - z/a_2} + \frac{A(a_3)}{1 - z/a_3},$$

where the  $a_i$ s are the roots of  $Q$  and  $A(z) = -1/Q'(z)$ . It follows by extracting coefficients that

$$h(n) = A(a_1)a_1^{-n} + A(a_2)a_2^{-n} + A(a_3)a_3^{-n}.$$

(Note that this is an actual equality valid for all  $n \geq 0$  and not an asymptotic result). Now, the roots of  $Q$  are approximately

$$a_1 = 0.5897545, \quad a_2 = -0.294877 - 0.872272i, \quad a_3 = -0.294877 + 0.872272i.$$

Since  $|a_2| = |a_3| = .9207 > |a_1|$ , it follows that the asymptotic behaviour is given by the term in  $a_1$ . We have proved the following theorem.



**Theorem 7** *The number  $h(n)$  of saturated stem-loops on  $[1, n]$  satisfies*

$$h(n) \sim 0.323954 \cdot 1.69562^n. \quad (32)$$

Convergence of the asymptotic limit in equation (32) is exponentially fast, so that when  $n = 20$ ,  $0.323954 \cdot 1.69562^n = 12504.2$ , while the exact number of saturated stem-loops on  $[1, 20]$  is  $h(20) = 12503$ .

### 3 Quasi-random saturated structures

In this section, we define a stochastic greedy process to generate *random* saturated structures, technically denoted *quasi-random saturated structures*. Our main result is that the expected number of base pairs in quasi-random saturated structures is  $0.0.340633 \cdot n$ , just slightly more than the expected number  $0.337361 \cdot n$  of all saturated structures. This suggests that the introduction of stochastic greedy algorithms and their asymptotic analysis may prove useful in other areas of random graph theory.

Consider the following stochastic process to generate a saturated structure. Suppose that  $n$  bases are arranged in sequential order on a line. Select the base pair  $(1, u)$  by choosing  $u$ , where  $\theta + 2 \leq u \leq n$ , at random with probability  $1/(n - \theta - 1)$ . The base pair joining 1 and  $u$

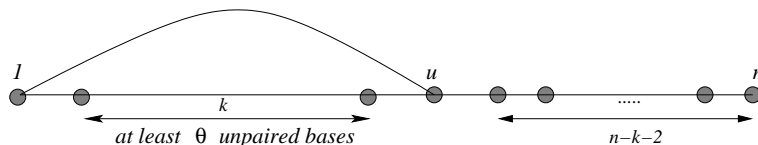


Figure 3: Base 1 is base-paired by selecting a random base  $u$  such there are at least  $\theta$  unpaired bases enclosed between 1 and  $u$ .

partitions the line into two parts. The left region has  $k$  bases strictly between 1 and  $u$ , where  $k \geq \theta$ , and the right region contains the remaining  $n - k - 2$  bases properly contained within endpoints  $k + 2$  and  $n$  (see Figure 3). Proceed recursively on each of the two parts. Observe that the secondary structures produced by our stochastic process will always base pair with the leftmost available base, and that the resulting structure is always saturated.

Before proceeding further, we note that the probability that the probability  $p_{i,j}$  that  $(i, j)$  is a base pair in a saturated structure is *not* the same as the probability  $q_{i,j}$  that  $(i, j)$  is a base pair in a quasi-random saturated structure. Indeed, if we consider saturated and quasi-random saturated structures on an RNA sequence of length  $n = 10$ , then clearly  $p_{1,5} = 1/29$  while clearly  $q_{1,5} = 1/8$ .<sup>¶</sup> Despite the very different base pairing probabilities when comparing saturated with quasi-random saturated structures, it is remarkable that the expected number of base pairs over saturated and quasi-random saturated structures is numerically so close.

<sup>¶</sup>The web supplement contains a Python program to compute the number of saturated structures on  $n$ . Clearly  $p_{1,5} = \frac{s_3 \cdot s_5}{s_{10}}$ , where  $s_k$  denotes the number of saturated structures on an RNA sequence of length  $k$ . A computation from a Python program (see web supplement) shows that  $s_3 = 1$ ,  $s_5 = 5$  and  $s_{10} = 145$ , hence  $p_{1,5} = 5/145 = 1/29$ .

Let  $U_n^\theta$  be the expected number of base pairs of the saturated secondary structure generated by this recursive procedure. In general, we have the following recursive equation

$$U_n^\theta = 1 + \frac{1}{n - \theta - 1} \sum_{k=\theta}^{n-2} (U_k^\theta + U_{n-k-2}^\theta), \quad n \geq \theta + 2, \quad (33)$$

with initial conditions

$$U_0^\theta = U_1^\theta = \dots = U_{\theta+1}^\theta = 0, \quad U_{\theta+2}^\theta = U_{\theta+3}^\theta = 1. \quad (34)$$

If we write equation (33) for  $U_{n+1}^\theta$  and substitute in it the value for  $U_n^\theta$  we derive

$$\begin{aligned} U_{n+1}^\theta &= 1 + \frac{1}{n - \theta} \sum_{k=\theta}^{n-1} (U_k^\theta + U_{n-k-1}^\theta) \\ &= 1 + \frac{1}{n - \theta} \left( U_{n-1}^\theta + U_{n-\theta-1}^\theta + \sum_{k=\theta}^{n-2} (U_k^\theta + U_{n-k-2}^\theta) \right) \\ &= 1 + \frac{1}{n - \theta} (U_{n-1}^\theta + U_{n-\theta-1}^\theta) + \frac{n - \theta - 1}{n - \theta} (U_n^\theta - 1). \end{aligned}$$

If we multiply out by  $n - \theta$  and simplify we obtain

$$(n - \theta)U_{n+1}^\theta = 1 + (n - \theta - 1)U_n^\theta + U_{n-1}^\theta + U_{n-\theta-1}^\theta, \quad (35)$$

which is valid for  $n \geq \theta + 1$ .

### 3.1 Asymptotic behavior

We now look at asymptotics. In particular we prove the following result.

**Theorem 8** *Let  $U_n^\theta$  denote the expected number of base pairs for quasi-random saturated structures of an RNA sequence of length  $n$ . Then for fixed  $\theta$ , and as  $n \rightarrow \infty$*

$$U_n^\theta \sim K_\theta \cdot n \quad \text{with} \quad K_\theta = e^{-1-H_{\theta+1}} \int_0^1 e^{t+(t+t^2/2+\dots+t^{\theta+1}/(\theta+1))} dt, \quad (36)$$

where  $H_{\theta+1} = 1 + \frac{1}{2} + \dots + \frac{1}{\theta+1}$  is the  $(\theta + 1)$ th harmonic number.

The first few values can easily be obtained numerically and we have

$$K_1 = 0.340633, \quad K_2 = 0.285497, \quad K_3 = 0.247908, \quad K_4 = 0.220308, \quad K_5 = 0.199018.$$

**Proof.** For fixed integer  $\theta$ , the recurrence (35) is linear with polynomial coefficients. It is a classical result that the generating functions of solutions of such recurrences satisfy linear differential equations. This is obtained by applying the following rules: if  $U(z) = \sum_{n \geq 0} u_n z^n$ , then

$$\sum_{n \geq 0} n u_n z^n = zU'(z), \quad \sum_{n \geq 0} u_{n+k} z^n = \frac{1}{z^k} (U(z) - u_0 - u_1 z - \dots - u_{k-1} z^{k-1}).$$

Starting from (35), we first shift the index by  $\theta + 1$  and apply these rules together with the initial conditions (34) to get

$$(n + \theta + 2)U_{n+\theta+2}^\theta - (\theta + 1)U_{n+\theta+2}^\theta = 1 + (n + \theta + 1)U_{n+\theta+1}^\theta - (\theta + 1)U_{n+\theta+1}^\theta + U_{n+\theta}^\theta + U_n^\theta,$$

$$\frac{1}{z^{\theta+2}}zy' - (\theta + 1)\frac{y}{z^{\theta+2}} = \frac{1}{1-z} + \frac{1}{z^{\theta+1}}zy' - (\theta + 1)\frac{y}{z^{\theta+1}} + \frac{y}{z^\theta} + y.$$

Finally, this simplifies to

$$z(1-z)y' + ((\theta + 1)(z-1) - z^2 - z^{\theta+2})y = \frac{z^{\theta+2}}{1-z}. \quad (37)$$

This is a first order non-homogeneous linear differential equation. The homogeneous part

$$z(1-z)W' + ((\theta + 1)(z-1) - z^2 - z^{\theta+2})W = 0$$

is solved by integrating a partial fraction decomposition

$$\frac{W'(z)}{W(z)} = \frac{\theta + 1}{z} - \frac{z}{z-1} - \frac{z^{\theta+1}}{z-1}$$

$$= \frac{\theta + 1}{z} + \frac{2}{z-1} - 1 - (1 + z + \dots + z^\theta)$$

$$\log W = (\theta + 1) \log z - 2 \log(1-z) - z - (z + z^2/2 + \dots + z^{\theta+1}/(\theta + 1)),$$

$$W(z) = \frac{z^{\theta+1}}{(1-z)^2} e^{-z-(z+z^2/2+\dots+z^{\theta+1}/(\theta+1))}.$$

From there, variation of the constant gives the following expression for the generating function:

$$y = \frac{z^{\theta+1}}{(1-z)^2} e^{-z-(z+z^2/2+\dots+z^{\theta+1}/(\theta+1))} \int_0^z e^{t+(t+t^2/2+\dots+t^{\theta+1}/(\theta+1))} dt.$$

Because the exponential is an entire function, we readily find that the only singularity is at  $z = 1$ , where  $y \sim K/(1-z)^2$  with  $K$  as in the statement of the theorem. The proof is completed by the use of Theorem 1.  $\blacksquare$

Note that the asymptotic expected number of base pairs in quasi-random saturated structures with  $\theta = 1$  is  $0.340633 \cdot n$ , while by Theorem 6 the asymptotic expected number of base pairs in saturated structures is  $0.337361 \cdot n$ , just very slightly less. This result points out that the stochastic greedy method performs reasonably well in sampling saturated structures, although the stochastic process tends not to sample certain (rare) saturated structures having a less than average number of base pairs.

The stochastic process used to construct quasi-random saturated structures iteratively base-pairs the leftmost position in each subinterval. One can imagine a more general stochastic method of constructing saturated structures, described as follows. Generate an initial list  $L$  of all allowable base pairs  $(i, j)$  with  $1 \leq i < j \leq n$  and  $j \geq i + \theta + 1$ . Create a saturated structure by repeatedly picking a base pair from  $L$ , adding it to an initially empty structure  $S$ , then removing from  $L$  all base pairs that form a crossing (pseudoknot) with the base pair just selected. This ensures that the next time a base pair from  $L$ , it can be added to  $S$  without violating the definition of secondary structure. Iterate this procedure until  $L$  is empty to form the stochastic saturated structure  $S$ .

Taking an average over 100 repetitions, we have computed the average number of base pairs and the standard deviation for  $n = 10, 100, 1000$ . Results are  $\mu = 0.323$ ,  $\sigma = 0.0604$  for  $n = 10$ ,  $\mu = 0.3526$ ,  $\sigma = 0.0386$  for  $n = 100$  and  $\mu = 0.35618$ ,  $\sigma = 0.0361$  for  $n = 1000$ . This clearly is a different stochastic process than that used for quasi-random saturated structures.

## 4 Conclusion

In this paper we applied the DSV methodology and the Flajolet-Odlyzko theorem to asymptotic enumeration problems concerning canonical and saturated secondary structures. For instance, we showed that the expected number of base pairs in canonical RNA secondary structures is equal to  $0.31724 \cdot n$ , which is far less than the expected number  $0.495917 \cdot n$  of base pairs over all secondary structures, the latter which follows from Theorem 4.19 of [12]. This may provide a theoretical explanation for the speed-up observed for Vienna RNA Package when restricted to canonical structures [1].

Additionally, we computed the asymptotic number  $1.07427 \cdot n^{-3/2} \cdot 2.35467^n$  of saturated structures, the expected number  $0.337361 \cdot n$  of base pairs of saturated structures and the asymptotic number  $0.323954 \cdot 1.69562^n$  of saturated stem-loop structures. We then considered a natural stochastic greedy process to generate quasi-random saturated structures, and showed surprisingly that the expected number of base pairs of is  $0.340633 \cdot n$ , a value very close to the expected number  $0.337361 \cdot n$  of base pairs of all saturated structures. Finally, we apply a theorem of Drmota [9] to show that the density of states for [all resp. canonical resp. saturated] secondary structures is asymptotically Gaussian. In the Appendix, we provide some background discussion and some technical developments that determine certain structural properties for quasi-random saturated RNA secondary structures – in particular stem length and number of external loops for the case when the associated tree graph is constructed using either the uniform or Zipf distributions.

## Acknowledgements

We would like to thank Yann Ponty, for suggesting that Drmota’s work can be used to prove that the density of states for secondary structures is Gaussian. Thanks as well to two anonymous referees, whose comments led to important improvements in this paper. Figure 2 is due to W.A. Lorenz, and first appeared in the joint article Lorenz et al. [17].

Funding for the research of P. Clote was generously provided by the Foundation Digiteo - Triangle de la Physique and the National Science Foundation DBI-0543506 and DMS-0817971. Additional support is gratefully acknowledged to the Deutscher Akademischer Austauschdienst for a visit to Martin Vingron’s group in the Max Planck Institute of Molecular Genetics. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation. Funding for the research of E. Kranakis was generously provided by the Natural Sciences and Engineering Research Council of Canada (NSERC) and Mathematics of Information Technology and Complex Systems (MITACS). Funding for the research of B. Salvy was provided by Microsoft Research-Inria Joint Centre.

## Appendix

In this appendix, we prove some technical results concerning expected stem length and the number of external loops of quasi-random saturated structures defined by a different stochastic process.

### Structural properties of random saturated secondary structures

Here we consider a different stochastic process in generating random saturated secondary structures, for which we determine expected number of external loops and expected stem size, formally defined below. Assume that  $n$  bases, numbered  $1, 2, \dots, n$  are arranged on a line, as depicted in Figure 3, and consider the stochastic process defined in Section 3 for generating random saturated secondary structures. Given secondary structure  $S$ , an *external base pair* is a base pair  $(i, j) \in S$ , which is not interior to any other base pair of  $S$ ; i.e. there is no  $(x, y) \in S$  with the property that  $x < i < j < y$ . A sequence of external base pairs is a sequence  $(a_i, b_i)$ ,  $i = 1, 2, \dots, k$  such that  $a_i < b_i < a_{i+1} < b_{i+1}$ , for all  $i < k$ , and for which each  $(a_i, b_i)$  is external. The base pairs  $(a_i, b_i)$  are said to *close* the corresponding *external loops*; see Figure 4. The *number of external loops* of a given secondary structure  $S$  is defined to be the maximum number of external base pairs in  $S$ . We define a *stem* of

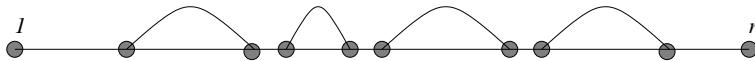


Figure 4: A sequence of external base pairs.

length  $k$  to be a sequence of nested base pairs (see Figure 5)  $(a_i, b_i)$ ,  $i = 1, 2, \dots, k$ , such that  $a_i < a_{i+1} < b_{i+1} < b_i$ , for all  $i < k$ . The *stem length* of a given secondary structure  $S$  is defined here to be the maximum length of all stems in  $S$ ; i.e. the maximum number of nested base pairs in  $S$ . Our study of structural properties of random saturated secondary structures

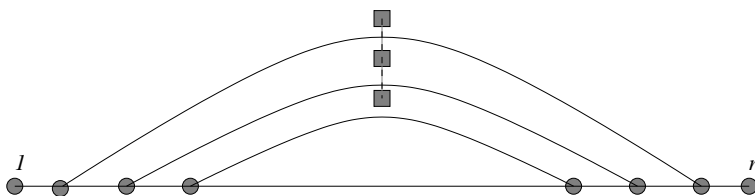


Figure 5: A sequence of nested base pairs.

is facilitated by defining a graph that resembles the graph on page 333 of [27]; however, note that the formal definition is slightly different than that of [27]. Given a secondary structure  $S$  on the nucleotide sequence  $[1, n]$ , define the associated graph  $G(S) = (V, E)$ , whose vertex set  $V$  consists of base pairs  $v = (i, j)$  in  $S$ , and whose undirected edge set  $E$  consists of pairs  $\{v, v'\}$  of nested vertices,  $v = (i, j)$  and  $v' = (i', j')$ , that can directly *see* each other; i.e.  $\{v, v'\} \in E$  exactly when  $i < i' < j' < j$  and there does not exist a base pair  $(x, y) \in S$ , such that  $i < x < i' < j' < y < j$ , or vice-versa with the roles of  $v, v'$  reversed. Figure 6 depicts the graph  $G(S)$  associated with the saturated secondary structure  $S$ . In general

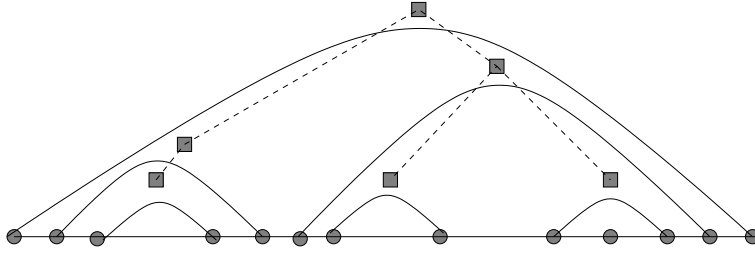


Figure 6: The tree associated with the given set of base pairs.

$G(S)$  is a forest; i.e., a set of trees. In the sequel we determine the size of several structural parameters of random saturated secondary structures, in particular, expected stem length and expected number of external loops. These parameters are studied both for the uniform and Zipf distributions. Before proceeding any further, we first define the probability distributions to be considered.

### Probability distributions

*Zipf's law* is the observation first made by the deceased Harvard linguist, George Kingsley Zipf, that the frequency  $p_i$  of English words, when graphed against their rank  $i$  (in the list of English words sorted in decreasing order with respect to frequency), obeys the power law  $p_i \approx i^{-\alpha}$ . More generally, Zipf's law is the statement of a power law, when plotting frequency against rank (Zipf's first law) or when plotting frequency against reverse rank (Zipf's second law). In bioinformatics, Zipf's law has been observed in the frequency/rank plot of differentially expressed gene in microarray data [16], as well as in the frequency/rank plot for protein structures [2], where there are few very frequent structures and very many rare structures are found. In the remainder of the paper, we consider probability distributions related to Zipf's law.

A node, say  $u$ , is chosen at random with the  $\alpha$ -Zipf distribution, if the probability that a given base pair  $(1, u)$  is chosen is equal to  $\frac{1}{(u-1)^\alpha H_\alpha(n-1)}$ , where

$$H_\alpha(n-1) = \sum_{k=1}^{n-1} \frac{1}{k^\alpha}$$

is defined to be the  $\alpha$ -harmonic number of  $n-1$ . As in equation (33), the expected number of base pairs, denoted by  $E_n^\theta$ , for random saturated secondary structures on  $n$  bases, generated by the  $\alpha$ -Zipf stochastic process, satisfies the following recursive formula

$$E_n^\theta = 1 + \frac{1}{H_\alpha(n-\theta-1)} \sum_{k=\theta}^{n-2} \frac{1}{(k-\theta+1)^\alpha} (E_k^\theta + E_{n-k-2}^\theta), \quad (38)$$

for all  $n \geq \theta + 2$ .

Observe that when  $\alpha = 0$  the  $\alpha$ -Zipf distribution is the same as the uniform distribution, while if  $\alpha = 1$ , we have the (classical) Zipf distribution [29]. Moreover, observe that as  $\alpha$  increases, "shorter" base pairs are being selected with higher probability by the stochastic process described in equation (38).

The stochastic process of generating random saturated secondary structures, according to equation (38), is of the “divide-and-conquer” type, very common in computer science, where well-known algorithms such as QUICKSORT choose a division point according to the uniform distribution. Stochastic algorithms of this kind have been intensively studied for the uniform distribution. Known results suggest that the probability distribution for the number of base pairs in random saturated structures, generated by the earlier described stochastic process (uniform choice of base pairs) is asymptotically Gaussian (see [6] and [13]). We also note that structural features of trees have been well studied including the expected depth and the exact distribution of the depth; see, for instance, [20, 5, 6]. In the sequel, we consider a random binary search tree with  $n$  nodes obtained by inserting  $n$  i.i.d. random variables  $X_1, \dots, X_n$ . Careful analysis of [6] and [5] implies our results in Section 4 for the uniform distribution. However we will use a different and simpler technique that enables the analysis not only for the uniform distribution in Subsection 4 but also for the Zipf distribution in Subsection 4.

An important observation concerns the threshold  $\theta$  considered above. All the results proved in this section are “upper bounds” and therefore it is easily seen that they are valid for any threshold  $\theta \geq 0$ . Therefore to simplify proofs in the sequel we consider the case of threshold  $\theta = 0$ .

## Uniform distribution

The main theorem of this section concerns stem length and number of external loops of random saturated structures  $S$ , generated by a natural stochastic process associated with the tree graph  $G(S)$ . Throughout the remainder of the paper, we state results in terms of *random saturated* structures, although we intend to mean only those structures generated by the stochastic process associated with the graph  $G(S)$ , although we will distinguish between the uniform and  $\alpha$ -Zipf variant of the stochastic process. Without this convention, statements of lemmas and theorems would be too cumbersome.

**Theorem 9** *With high probability, the number of external loops and the maximum stem length of random saturated structures generated by the uniform distribution variant is  $O(\log n)$ .*

**Proof.** Before we give the proof of the main theorem it will be necessary to give the proof of two lemmas. In the first lemma we consider the maximum number of external loops.

**Lemma 1** *With high probability, the maximum number of external loops is  $O(\log n)$ .*

**Proof.** We define a sequence of random variables  $X_1, X_2, \dots, X_t$  by induction as follows. Let  $X_1$  be the random variable selecting a base  $k$  chosen among  $2, 3, \dots, n$  randomly and independently with the uniform distribution in order to form a base pair  $(1, k)$ . By induction, assume that  $X_1, \dots, X_t$  have been defined. Let  $X_{t+1}$  be the random variable selecting a base  $k$  chosen among  $X_t + 2, X_t + 3, \dots, n$  randomly and independently with the uniform distribution in order to form a base pair  $(X_t + 1, k)$ . Next we compute  $\mathbb{E}[X_t]$ , for all  $t$ . Indeed, observe that  $\mathbb{P}[X_1 = k] = \frac{1}{n-2}$  and

$$\mathbb{E}[X_1] = \sum_{i=2}^n i \cdot \frac{1}{n-1} = \frac{1}{n-1} \sum_{i=2}^n i = \frac{1}{n-1} \left( \frac{n(n+1)}{2} - 1 \right)$$

Next we compute the conditional probability

$$\begin{aligned}\mathbb{E}[X_{t+1}|X_t = k] &= \sum_{i=k+2}^n i \cdot \mathbb{P}[X_{t+1} = i|X_t = k] = \sum_{i=k+2}^{n-1} i \cdot \frac{1}{n-k-2} \\ &= \frac{1}{n-k-2} \sum_{i=k+2}^{n-1} i = \frac{1}{n-k-2} \left( \sum_{i=0}^{n-1} i - \sum_{i=0}^{k+1} i \right) = \frac{n+k+1}{2}.\end{aligned}$$

Finally, we can calculate

$$\begin{aligned}\mathbb{E}[X_{t+1}] &= \mathbb{E}[\mathbb{E}[X_{t+1}|X_t]] = \sum_k \mathbb{E}[X_{t+1}|X_t = k] \cdot \mathbb{P}[X_t = k] = \sum_k \frac{n+k+1}{2} \cdot \mathbb{P}[X_t = k] \\ &= \frac{n+1}{2} + \frac{1}{2} \sum_k k \cdot \mathbb{P}[X_t = k] = \frac{n+1}{2} + \frac{1}{2} \mathbb{E}[X_t] = \frac{n+1}{2} \cdot (1 + 2^{-1} + \dots + 2^{-t}) \\ &= (n+1) (1 - 2^{-t-1}).\end{aligned}$$

We are interested in determining the behavior of the random variable, whose value is the number of external loops in random saturated structures.

$$T_n = \min\{t : X_{t+1} \geq n+1\}. \quad (39)$$

From this we derive

$$\begin{aligned}\mathbb{P}[T_n > t] &= \mathbb{P}[X_{t+1} < n+1] = \mathbb{P}[n+1 - X_{t+1} > 0] \leq \mathbb{E}[n+1 - X_{t+1}] \\ &= n+1 - \mathbb{E}[X_{t+1}] = n+1 - (n+1) (1 - 2^{-t-1}) = (n+1)2^{-t-1}.\end{aligned}$$

In particular,  $\mathbb{P}[T_n > (1 + \epsilon) \log n] \leq n^{-\epsilon}$ . This completes the proof of Lemma 1. ■

Next we prove the following lemma.

**Lemma 2** *With high probability, the maximum stem length is  $O(\log n)$ .*

**Proof.** According to the recursive construction, at each stage after a base pair is chosen at random in the subsequent stages, base pairs are nested within this base pair. Therefore, the maximum stem length equals the maximum number of nested base pairs. This latter number can also be obtained as follows. We define the following sequence  $Y_1, Y_2, \dots, Y_t$  of random variables. A base is chosen among  $2, 3, \dots, n$  randomly and independently with the uniform distribution. Let  $Y_1$  be the resulting random variable. By induction, assume that  $Y_1, \dots, Y_t$  have been defined. To define the random variable  $Y_{t+1}$ , a base is chosen among  $t+2, \dots, Y_t - 1$  randomly and independently with the uniform distribution. Clearly, this procedure halts when  $Y_t = 1$  and it follows that the maximum number of nested base pairs is also the number  $t$  of iterations before halting. Therefore we are interested in knowing the behavior of the random variable

$$T'_n = \min\{t : Y_t > 0\}. \quad (40)$$

Observe that since by definition  $Y_{i+1}$  is chosen among  $i+2, i+3, \dots, Y_i - 1$  randomly and independently with the uniform distribution, for any integer  $k \geq i+2$ ,  $\mathbb{E}[Y_{i+1}|Y_i = k] = \frac{k-i-1}{2}$ .



Consider the random variable  $\mathbb{E}[Y_{i+1}|Y_i]$  whose value at  $k$  is equal to  $\mathbb{E}[Y_{i+1}|Y_i = k]$ . Using well-known identities on conditional probabilities, we can derive the following equalities.

$$\begin{aligned}
\mathbb{E}[Y_{i+1}] &= \mathbb{E}[\mathbb{E}[Y_{i+1}|Y_i]] \\
&= \sum_k \mathbb{E}[Y_{i+1}|Y_i = k] \cdot \mathbb{P}[Y_i = k] \\
&= \sum_k \frac{k-i-1}{2} \cdot \mathbb{P}[Y_i = k] \\
&= \frac{1}{2} \sum_k k \cdot \mathbb{P}[Y_i = k] - \frac{i+1}{2} \\
&= \frac{1}{2} \mathbb{E}[Y_i] - \frac{i+1}{2}.
\end{aligned}$$

In particular, since  $E[Y_1] = \frac{n-2}{2}$ , we conclude that  $\mathbb{E}[Y_t] \leq (1/2)^t \cdot n$ . Finally, we can derive  $\mathbb{P}[T'_n > t] = \mathbb{P}[Y_t > 0] \leq \mathbb{E}[Y_t] \leq (1/2)^t \cdot n$ . It follows that  $\mathbb{P}[T'_n > (1 + \epsilon) \log n] \leq n^{-\epsilon}$ .

We are not yet completely done with the proof of Lemma 2. The proof shows that with high probability, the leftmost sequence of base pairs given by the recursive construction has length at most  $O(\log n)$ . We would like to prove the same for any sequence of nested base pairs. To this effect, define random intervals  $I_s$ , where  $s$  is a finite sequence of 0s and 1s, by induction on the length of  $s$ . Consider the interval  $I_\emptyset = [1, n]$ . Assuming that  $I_s = [a_s, b_s]$  has already been defined, we consider a random process that splits it at random into two subintervals, i.e., choose an integer  $r \in I_s$  randomly and independently with the uniform distribution and let  $I_{s0} = [a_s, r]$  and  $I_{s1} = [r+1, b_s]$ . Since  $E[|I_{sb}|] \leq \frac{1}{2} \cdot E[|I_s|]$  it follows that the expected length of  $I_s$  is at most  $2^{-|s|}$ . Now consider the random variable  $T''_n$  which is defined as follows  $T''_n = \min\{k : \exists s(|s| = k \ \& \ I_s = \emptyset)\}$  and observe that  $T''_n > k$  if and only if  $\forall s(|s| = k \Rightarrow I_s \neq \emptyset)$ . Therefore

$$\begin{aligned}
\mathbb{P}[T''_n > k] &= \mathbb{P}[\min_{k:|s|=k} |I_s| > 0] \leq \mathbb{E}[\min_{k:|s|=k} |I_s|] \\
&\leq \mathbb{E}[|I_s|], \text{ (for all sequences } s \text{ such that } |s| = k) \\
&\leq 2^{-k}.
\end{aligned}$$

As a consequence we conclude that  $\mathbb{P}[T''_n > (1 + \epsilon) \log n] \leq n^{-\epsilon}$ . This completes the proof of Lemma 2. ■

Finally, we can complete the proof of the main result of Theorem 9 since this is now immediate from Lemmas 1 and 2. ■

## Zipf distribution

It is possible to consider other probability distributions like Zipf and generalized  $a$ -Zipf. The Zipf distribution (first considered in [29]) is perhaps the most interesting because it favors base pairs at a shorter distance. A base pair  $(1, u)$ , is chosen at random with the Zipf distribution. I.e., the probability that the base pair  $(1, u)$  is selected is equal to  $\frac{1}{(u-1)H(n-1)}$ , where

$$H(n-1) = \sum_{k=1}^{n-1} \frac{1}{k}$$

is defined to be the  $(n-1)$ st harmonic number. As before, the chord joining 1 and  $u$  partitions the ring into two parts. One part has  $k$  bases between 1 and  $u$ , where  $k \leq n-2$ , and the other part has the remaining  $n-k-2$  bases (see Figure 3).

Define  $Z_n$  to be the expected number of base pairs of a random saturated secondary structure with  $n$  bases, where  $n \geq 2$ . A base pair  $(1, u)$  is added as follows. Select  $u \geq 2$  at random among  $2, 3, \dots, u-1$  with probability  $\frac{1}{(u-1)H(n-1)}$ . This gives rise to the following formula

$$Z_n = 1 + \frac{1}{H(n-1)} \sum_{k=0}^{n-2} \frac{1}{k+1} (Z_k + Z_{n-k-2}), \quad (41)$$

for all  $n \geq 2$ . The main theorem of this section concerns the overall structure of random secondary structures.

**Theorem 10** *With high probability, random saturated secondary structure generated by the Zipf distribution have  $O(\log^2 n)$  external loops and stem length  $O(\log n / \log \log n)$ .*

**Proof.** Before we give the proof, it will be necessary to give the proof of two lemmas. In the first lemma we look at the number of external loops.

**Lemma 3** *With high probability, the number of external loops is  $O(\log^2 n)$ .*

**Proof.** We define a sequence of random variables  $X_1, X_2, \dots, X_t$  by induction as follows. Let  $X_1$  be the random variable resulting when the base pair  $(1, k)$  is formed by a selecting a base  $k$  among  $2, 3, \dots, n-1$  randomly and independently with the uniform distribution. By induction, assume that  $X_1, \dots, X_t$  have been defined. Let  $X_{t+1}$  be the random variable resulting when the base pair  $(X_t + 1, k)$  is formed by selecting a base  $k$  is chosen among  $X_t + 2, X_t + 3, \dots, n-1$  randomly and independently with the uniform distribution. Next we compute  $\mathbb{E}[X_t]$ , for all  $t$ . Indeed, observe that  $\mathbb{P}[X_1 = k] = \frac{1}{(k-1)H(n-2)}$  and

$$\mathbb{E}[X_1] = \sum_{i=2}^{n-1} (i-1) \cdot \frac{1}{(i-1)H(n-2)} = \frac{n-2}{H(n-2)}.$$

Next we compute the conditional probability

$$\begin{aligned} \mathbb{E}[X_{t+1} | X_t = k] &= \sum_{i=k+2}^{n-1} i \cdot \mathbb{P}[X_{t+1} = i | X_t = k] \\ &= \sum_{i=k+2}^{n-1} i \cdot \frac{1}{(i-k-1)H(n-k-2)} \\ &= \frac{1}{H(n-k-2)} \sum_{i=k+2}^{n-1} \frac{i}{i-k-1} \\ &= \frac{1}{H(n-k-2)} \sum_{i=k+2}^{n-1} \left( \frac{i-k-1}{i-k-1} + \frac{k+1}{i-k-1} \right) \\ &= \frac{n-k-2}{H(n-k-2)} + (k+1). \end{aligned}$$

Finally, we can calculate

$$\begin{aligned}
\mathbb{E}[X_{t+1}] &= \mathbb{E}[\mathbb{E}[X_{t+1}|X_t]] \\
&= \sum_k \mathbb{E}[\mathbb{E}[X_{t+1}|X_t = k]] \cdot \mathbb{P}[X_t = k] \\
&= \sum_k \left( (k+1) + \frac{n-k-2}{H(n-k-2)} \right) \cdot \mathbb{P}[X_t = k] \\
&= 1 + \mathbb{E}[X_t] + \sum_k \frac{n-k-2}{H(n-k-2)} \cdot \mathbb{P}[X_t = k] \\
&\geq 1 + \mathbb{E}[X_t] + \frac{1}{H(n-2)} \sum_k (n-k-2) \cdot \mathbb{P}[X_t = k] \\
&= 1 + \mathbb{E}[X_t] + \frac{1}{H(n-2)} (n-2 - \mathbb{E}[X_t]) \\
&\geq \frac{n-2}{H(n-2)} + \left( 1 - \frac{1}{H(n-2)} \right) \mathbb{E}[X_t].
\end{aligned}$$

Elementary calculations using this last inequality show that

$$\mathbb{E}[X_{t+1}] \geq (n-2) \left( 1 - \left( 1 - \frac{1}{H(n-2)} \right)^{t+2} \right)$$

We are interested in determining the behavior of the random variable, whose value is the number of external loops; i.e. the size of the largest sequence of external base pairs.

$$T_n = \min\{t : X_{t+1} \geq n-2\}. \quad (42)$$

From this we derive

$$\begin{aligned}
\mathbb{P}[T_n > t] &= \mathbb{P}[X_{t+1} < n-2] = \mathbb{P}[n-2 - X_{t+1} > 0] \leq \mathbb{E}[n-2 - X_{t+1}] = n-2 - \mathbb{E}[X_{t+1}] \\
&\leq n-2 - (n-2) \left( 1 - \left( 1 - \frac{1}{H(n-2)} \right)^{t+2} \right) = (n-2) \left( 1 - \frac{1}{H(n-2)} \right)^{t+2}.
\end{aligned}$$

In particular, since  $H(n-2) \approx \ln n$  we conclude that  $\mathbb{P}[T_n > \epsilon \ln^2 n] \leq n^{-\epsilon}$ . This completes the proof of Lemma 3.  $\blacksquare$

The next result concerns the maximum stem length. We can prove the following result.

**Lemma 4** *With high probability, the maximum stem length is  $O(\log n / \log \log n)$ .*

**Proof.** According to the recursive construction, at each stage after a base pair is chosen at random in the subsequent stages base pairs are nested within this base pair. Therefore, the maximum stem length is equal to the maximum number of nested base pairs. This latter number can also be obtained, by investigating a sequence of random variables  $Y_1, Y_2, \dots, Y_t$ , defined as follows. Choose a base among  $2, 3, \dots, n-1$  randomly and independently with the uniform distribution. Let  $Y_1$  be the resulting random variable. By induction, assume that  $Y_1, \dots, Y_t$  have been defined. To define the random variable  $Y_{t+1}$ , a base is chosen among  $t+2, t+3, \dots, Y_t-1$  randomly and independently with the uniform distribution. Clearly,

this procedure halts when  $Y_t = 1$  and it follows that the maximum number of nested base pairs is also the number  $t$  of iterations before halting. Therefore we are interested to know the behavior of the random variable

$$T'_n = \min\{t : Y_t > 0\}. \quad (43)$$

Observe that since by definition  $Y_{i+1}$  is chosen among  $i+2, i+3, \dots, Y_i-1$  randomly and independently with the uniform distribution, for any integer  $k \geq i+2$ ,

$$\mathbb{E}[Y_{i+1}|Y_i = k] = \frac{k-i-1}{H(k-i-1)}.$$

Consider the random variable  $\mathbb{E}[Y_{i+1}|Y_i]$  whose value at  $k$  is equal to  $\mathbb{E}[Y_{i+1}|Y_i = k]$ . Using well-known identities on conditional probabilities we can derive the following inequalities.

$$\begin{aligned} \mathbb{E}[Y_{i+1}] &= \mathbb{E}[\mathbb{E}[Y_{i+1}|Y_i]] \\ &= \sum_k \mathbb{E}[\mathbb{E}[Y_{i+1}|Y_i = k]] \cdot \mathbb{P}[Y_i = k] \\ &= \sum_{k \geq i+2} \frac{k-i-1}{H(k-i-1)} \cdot \mathbb{P}[Y_i = k] \\ &\leq \sum_{k \geq i+2} \frac{k}{H(k)} \cdot \mathbb{P}[Y_i = k] \\ &\leq \frac{1}{H(i+2)} \mathbb{E}[Y_i], \end{aligned}$$

where we used the fact that the fraction  $n/H(n)$  is monotone increasing in  $n$ . In particular, since  $\mathbb{E}[Y_1] = \frac{n-2}{H(n-2)}$ , we conclude that  $\mathbb{E}[Y_t] \leq \frac{n-2}{H(t+1) \cdot H(t) \cdots H(2)}$ . Finally, we can derive

$$\mathbb{P}[T'_n > t] = \mathbb{P}[Y_t > 0] \leq \mathbb{E}[Y_t] \leq \frac{n-2}{H(t+1) \cdot H(t) \cdots H(2)} \leq \frac{n-2}{H(t/2)^{t/2}}.$$

In particular,

$$\mathbb{P}\left[T'_n > (1+\epsilon) \frac{\log n}{\ln \ln n}\right] \leq n^{-\epsilon}.$$

The proof shows that the leftmost sequence of base pairs given by the recursive construction of the random secondary structure has length at most  $O(\log n / \log \log n)$  with high probability. We would like to prove the same for any sequence of nested base pairs. A proof similar to the one presented above should work. This completes the proof of Lemma 4. ■

If we now combine Lemmas 3 and 4 we derive the proof of Theorem 10. ■

## References

- [1] A. F. Bompfunewerer, R. Backofen, S. H. Bernhart, J. Hertel, I. L. Hofacker, P. F. Stadler, and S. Will. Variations on RNA folding and alignment: lessons from Benasque. *J. Math. Biol.*, 56(1-2):129–144, January 2008.
- [2] E. Bornberg-Bauer. How are model protein structures distributed in sequence space? *Biophys. J.*, 73(5):2393–2403, November 1997.

- [3] N. Chomsky and M. P. Schützenberger. The algebraic theory of context-free languages. In P. Braffort and D. Hirschberg, editors, *Computer Programming and Formal Languages*, pages 118–161. North Holland, 1963.
- [4] P. Clote. Combinatorics of saturated secondary structures of RNA. *J. Comput. Biol.*, 13(9):1640–1657, November 2006.
- [5] L. Devroye. Universal limit laws for depths in random trees. *SIAM Journal on Computing*, 28(2):409–432, 1998.
- [6] L. Devroye. Limit laws for sums of functions of subtrees of random binary search trees. *SIAM Journal on Computing*, 32(1):152–171, 2003.
- [7] R. Donaghey and L. W Shapiro. Motzkin numbers. *J. Combin. Theory*, 23:291–301, 1977.
- [8] Michael Drmota. Asymptotic distributions and a multivariate Darboux method in enumeration problems. *Journal of Combinatorial Theory, Series A*, 67(2):169–184, 1994.
- [9] Michael Drmota. Systems of functional equations. *Random Structures and Algorithms*, 10:103–124, 1999.
- [10] P. Flajolet and A. M. Odlyzko. Singularity analysis of generating functions. *SIAM Journal of Discrete Mathematics*, 3:216–240, 1990.
- [11] I.L. Hofacker. Vienna RNA secondary structure server. *Nucleic Acids Res.*, 31:3429–3431, 2003.
- [12] I.L. Hofacker, P. Schuster, and P. Stadler. Combinatorics of RNA secondary structures. *Discr. Appl. Math.*, 88:207–237, 1998.
- [13] H.-K. Hwang and R. Neininger. Phase change of limit laws in the quicksort recurrence under varying toll functions. *SIAM Journal on Computing*, 31(6):1687–1722, 2002.
- [14] S. Lang. *Algebra*. Springer Verlage, 2002. Revised 3rd edition.
- [15] H.R. Lewis and C.H. Papadimitriou. *Elements of the Theory of Computation*. Prentice-Hall, 1997. Second edition.
- [16] W. Li and Y. Yang. Zipf’s law in importance of genes for cancer classification using microarray data. *J. theor. Biol.*, 219(4):539–551, December 2002.
- [17] W.A. Lorenz, Y. Ponty, and P. Clote. Asymptotics of rna shapes. *J Compu Biol.*, 2007. in press.
- [18] A. Meir and J. W. Moon. On an asymptotic method in enumeration. *Journal of Combinatorial Theory, Series A*, 51(1):77–89, 1989.
- [19] A. Meir and J.W. Moon. On an asymptotic method in enumeration. *Journal of Combinatorial Theory*, 51:77–89, 1989. Series A.
- [20] M. E. Nebel. Investigation of the Bernoulli model for RNA secondary structures. *Bull. Math. Biol.*, 66(5):925–964, September 2004.

- [21] R. Nussinov and A. B. Jacobson. Fast algorithm for predicting the secondary structure of single stranded RNA. *Proceedings of the National Academy of Sciences, USA*, 77(11):6309–6313, 1980.
- [22] A.M. Odlyzko. Asymptotic enumeration methods. In R.L. Graham D.E. Knuth O. Patashnik, editor, *Concrete Mathematics - A Foundation for Computer Science*, pages 1063–1230. Addison-Wesley, 1989.
- [23] R.L. Graham D.E. Knuth O. Patashnik. *Concrete Mathematics - A Foundation for Computer Science*. Addison-Wesley, 1989.
- [24] P. Flajolet R. Sedgewick. *Analytic Combinatorics*. Cambridge University, 2009. ISBN-13: 9780521898065.
- [25] P. R. Stein and M. S. Waterman. On some new sequences generalizing the Catalan and Motzkin numbers. *Discrete Mathematics*, 26:261–272, 1978.
- [26] M. Szymanski, M. Z. Barciszewska, J. Barciszewski, and V. A. Erdmann. 5S ribosomal RNA database Y2K. *Nucleic. Acids. Res.*, 28(1):166–167, January 2000.
- [27] M.S. Waterman. *Introduction to Computational Biology: Maps, sequences and genomes*. Chapman & Hall – CRC Press, 1995.
- [28] K. C. Wiese, E. Glen, and A. Vasudevan. JViz.Rna—a Java tool for RNA secondary structure visualization. *IEEE. Trans. Nanobioscience.*, 4(3):212–218, September 2005.
- [29] G. K. Zipf. *Human Behavior and the Principle of Least Effort*. Addison Wesley, 1949.
- [30] M. Zuker. RNA folding prediction: The continued need for interaction between biologists and mathematicians. In *Lectures on Mathematics in the Life Sciences*, volume 17, pages 87–124. Springer-Verlage, 1986.
- [31] M. Zuker. Mfold web server for nucleic acid folding and hybridization prediction. *Nucleic Acids Res.*, 31(13):3406–3415, 2003.