



HAL
open science

Endhered patterns in matchings and RNA

Célia Biane, Greg Hampikian, Sergey Kirgizov, Khaydar Nurligareev

► **To cite this version:**

Célia Biane, Greg Hampikian, Sergey Kirgizov, Khaydar Nurligareev. Endhered patterns in matchings and RNA. 2024. hal-04563757v1

HAL Id: hal-04563757

<https://hal.science/hal-04563757v1>

Preprint submitted on 21 Oct 2024 (v1), last revised 12 Nov 2024 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Endhered patterns in matchings and RNA

Célia Biane¹, Greg Hampikian², Sergey Kirgizov¹, Khaydar Nurligareev¹

¹LIB, Université de Bourgogne, Dijon, France

²CompGenomics, Box 1454, Boise, Idaho 83701, USA

Corresponding author: sergey.kirgizov@u-bourgogne.fr

October 16, 2024

Abstract

An *endhered (end-adhered) pattern* is a subset of arcs in matchings, such that the corresponding starting points are consecutive and the same holds for the ending points. Such patterns are in one-to-one correspondence with the permutations. We focus on the occurrence frequency of such patterns in matchings and native (real-world) RNA structures with pseudoknots. We present combinatorial results related to the distribution and asymptotic behavior of the pattern 21, which corresponds to two consecutive base pairs frequently encountered in RNA, and the pattern 12, representing the archetypal minimal pseudoknot. We show that in matchings these two patterns are equidistributed, which is quite different from what we can find in native RNAs. We also examine the distribution of endhered patterns of size 3, showing how the patterns change under the transformation called *endhered twist*. Finally, we compute the distributions of endhered patterns of size 2 and 3 in native secondary RNA structures with pseudoknots and discuss possible outcomes of our study.

1 Introduction

Ribonucleic acids (RNAs) are macromolecules fulfilling many biological functions: they code for protein, are involved in the regulation of gene expression, can have catalytic activities and store the genetic information of certain viruses. The structure of RNAs is defined at the primary level as sequences of four nucleotides: Adenine (A), Uracil (U), Guanine (G), and Cytosine (C). The secondary structure abstracts from the nature of the nucleotide and considers only the bonds forming between nucleotides during the synthesis of RNAs and shaping how the molecules folds in space. Two types of bonds are formed during the RNA folding process: phospho-diester bonds (known as strong bonds) are formed between pairs of consecutive nucleotides in the sequence forming the RNA chain, and hydrogen bonds (also known as weak bonds) are formed between pairs of nucleotides distant in the sequence. The secondary structure represents an intermediate level between the primary sequence and the shape, and has the advantage of being both relevant from a biological perspective and tractable from a computational point of view.

1.1 Models of RNA secondary structures

RNA secondary structures have been formalized as graphs primarily by Waterman [46] to tackle the problem of prediction of secondary structures from the primary structure, which is more easily measurable. Ponty [38] gives a variant of Waterman definition of RNA secondary structure without pseudoknots and having a minimal number θ of unpaired positions between pair positions. Formally, *Waterman-Ponty RNA secondary structure* S of size n is defined as a set of base-pairs $(i, j), 1 \leq i < j \leq n$, such that:

1. Each position is monogamous, $\forall (i, j) \neq (i', j') \in S : \{i, j\} \cap \{i', j'\} = \emptyset$.
2. Minimal distance θ between paired nucleotides, $\forall (i, j) \in S : j - i \geq \theta$.
3. No pseudoknot allowed, $\forall (i, j), (i', j') \in S, i < i' : (j' < j) \text{ or } (j < i')$.

These structures can be represented using the *dot-bracket notation*. A secondary structure of an n -nucleotide RNA is encoded as an n -length sequence of parentheses $\{“(”;$ ”)” $\}$ and dots “.”, where an open parenthesis represents a nucleotide paired to another nucleotide represented by a closed parenthesis, and dots correspond to unpaired nucleotides. *The extended dot-bracket notation* includes also other types of parentheses: “[”]”, “{”}”, “<”>”, “aA”, etc. The extended dot-bracket notation enables the representation of pairings in pseudoknots. Figure 1 shows 4 different representations of an example of RNA secondary structure.

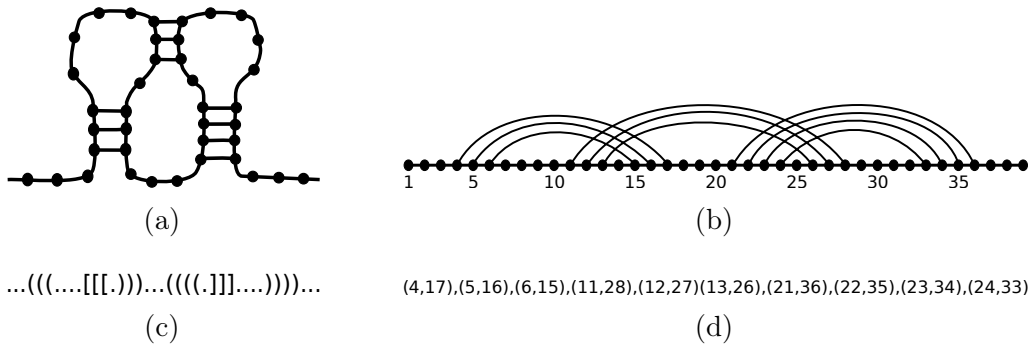


Figure 1: A drawing (a), an arc diagram (b), an extended dot bracket notation (c), and a set of pairs (d) representing an example of an RNA secondary structure.

Other models of RNA secondary structures were studied from combinatorial point of view. Haslinger and Standler [24] examined enumerative and asymptotic properties of *bi-secondary structures*, i.e., arc diagrams with arcs both in upper and lower part of the plane but without arc intersections. The concept of RNA shape was introduced by Giegerich, Voss, and Rehmsmeier [17] who developed an algorithm for the computation of minimum free energy RNA shapes. The number and asymptotics of RNA shapes was studied by Lorenz, Ponty, and Clote [30]. Reidys and Wang [40] considered a generalization of RNA shapes defined on k mutually-crossing arcs.

In this paper, we look at RNA structures with no restrictions on the number of arc crossings.

1.2 Notion of pattern in RNA secondary structures

At the biological level, common patterns have been observed in orthodox structures: single strand regions, hairpin and internal loops, bulges and various computational tools exist for detecting these patterns from primary sequence, including the work of Macke *et al.* [33]. More and more information is being gathered from three dimensional reconstructions of RNA molecules paving the way to a better comprehension of the laws governing the RNA folding process and the formation of RNA patterns.

Rødland [41] proposed a classification of RNA secondary structures in four level of abstractions: nucleotide, ladder, stem and collapsed level, based on the considered internal patterns. In his work, the nucleotide level corresponded to structures with arc diagrams containing unpaired nucleotides, the ladder level corresponded to structures abstracted from unpaired nucleotides, the stem level was abstracted from bulges and internal loops, and the collapsed level was abstracted from nested loops. Rødland studied different kind of pseudoknot patterns of increasing complexity: H-pseudoknot, double hairpin pseudoknot and pseudotrefoil. He counted these patterns in RNA secondary structures of increasing complexity and studied their asymptotics. He also showed that the theoretical number of pseudoknots in secondary structures is higher than in real secondary structure of the Rfam [21, 22] and PseudoBase [45] databases. Note that Rødland’s collapsed structures correspond to RNA shapes studied by Giegerich, Voss, and Rehmsmeier [17].

Quadrini [39] addressed the problem of searching a given structural pattern, defined as a sequence of crossing loops in a RNA secondary structure or shape and characterized by arbitrary number of pseudoknots. She proposed polynomial time algorithms for their identification. A paper by Gan, Pasquali, and Schlick [16] studied RNA structures and their patterns using graph-based representations. In the future, it would be interesting to compare the relationships between different kinds of patterns.

In this paper, we study the formation of endhered patterns (formally defined below) in matchings. In Section 2, we do it from a theoretical perspective. In Section 3, we observe the number of occurrences of such patterns in RNA secondary structures derived from experimentally determined 3D RNA structures. We conclude by discussing possible outcomes of our study in Section 4.

2 Endhered patterns in matchings

2.1 Basic definitions

By (perfect) *matching* of size n , we mean the sequence of $2n$ points $(1, 2, \dots, 2n - 1, 2n)$ endowed with a set of n arcs, such that every point is linked to one and only one another point. Figure 2a shows examples of matchings of small sizes. Matchings of size n can be considered as fixed-point-free involutions in the symmetric group S_{2n} . Thus, the matching with 4 arcs at the bottom of Figure 2a can be represented by the permutation 36154287 , which is the product of disjoint transpositions $(13)(26)(45)(78)$.

For a positive integer n , any matching of size n can be uniquely constructed from some matching of size $n - 1$ by the following procedure. We add a new arc starting at the left of the already existing $2(n - 1)$ points and ending at some of $2(n - 1) + 1$ possible new

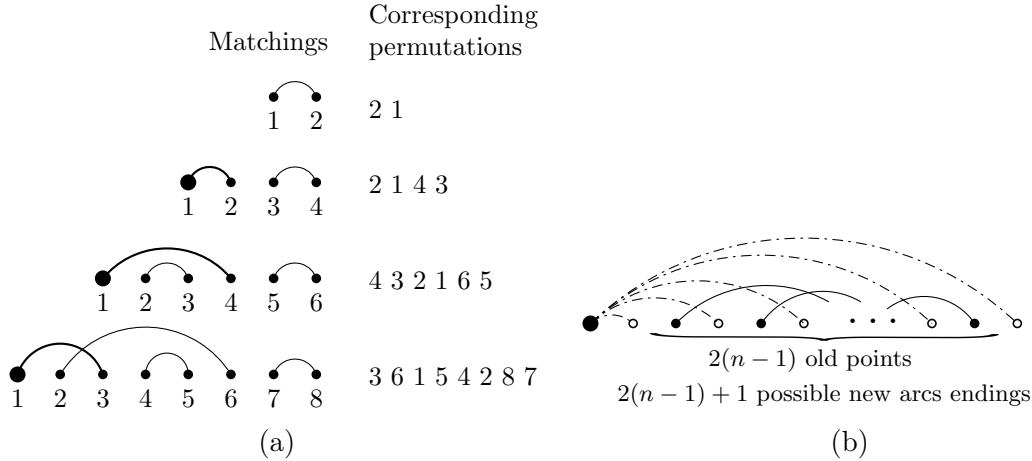


Figure 2: An example of matching construction, corresponding permutations (a), and a schema of recursive construction of matchings (b).

positions (see Figure 2b). This observation leads us to the following recurrence relation for the number of matchings of size n : $a_n = (2n - 1)a_{n-1}$ with $a_0 = 1$. As a consequence, $a_n = (2n - 1)!! = (2n - 1) \cdot (2n - 3) \cdot \dots \cdot 3 \cdot 1$. The corresponding sequence starting with 1, 1, 3, 15, 105, 945, 10395 is known as [A1147](#) in Sloane's Encyclopedia [35].

In mathematical literature, matchings often appear in different contexts, from the representation theory of Lie algebras [8] to the geometry of moduli spaces of flat connections on surfaces [2]. Efficient generating algorithms for involutions with (without) fixed points were established by Vajnovszki [44]. Ardnt's book [3] also presents interesting generating algorithms for involutions and other combinatorial structures.

Different kinds of patterns in matchings are being actively studied in combinatorics for the past twenty years. Initially, the interest to this topic came from the rapidly developing study of permutation patterns, since a matching can be thought as a permutation of a specific form. From this perspective, we say that a matching σ is a *pattern* in a matching μ if σ can be obtained from μ by deleting some of its arcs (and consistently relabelling the remaining vertices). For instance, Chen, Deng, Du, Stanley, and Yan [10] studied distributions of crossings and nestings. Jelínek [26], as well as Bloom and Elizalde [7], considered pattern avoiding matchings in the case when σ is a permutational matching of size 3. An extension for more general patterns was elaborated by Cervetti and Ferrari [9], while other authors, such as Chen, Mansour and Yan [11], Jelínek and Mansour [27], considered partial patterns.

As we see, what matters in the above investigations is the relative positions of arcs that form a pattern. At the same time, the distances between starting and ending points of these arcs are not fixed. The main object of our study concerns specific restrictions imposed on the arcs. Namely, the starting points of a pattern, as well as its ending points, form an interval, while the distance between these two intervals may vary. We call such patterns *endhered* (end-adhered) to emphasize the nature of these restrictions.

Definition 2.1. An *endhered pattern* is a matching, such that the starting point of any of its arcs precedes the ending point of any other arc. In other words, a matching of size p written as a permutation $\sigma = \sigma_1 \dots \sigma_{2p}$ is an endhered pattern if $\pi = \sigma_{p+1} \dots \sigma_{2p}$ is a

permutation of size p (such matchings are also called *permutational*). Figure 3a presents an example of an endhered pattern of size 3. In the following, we identify endhered patterns with the corresponding permutations.

We say that a matching $\mu = \mu_1 \dots \mu_{2n}$ contains an endhered pattern $\pi = \pi_1 \dots \pi_p$ at position $(i + 1, j + 1)$, where $i \geq 0$ and $i + p \leq j \leq n - p$, if

$$\mu_{s+i} = \pi_s^{-1} + j, \quad s = 1, \dots, p$$

(here, $\pi^{-1} = \pi_1^{-1} \dots \pi_p^{-1}$ is the inverse to the permutation π). In other words, μ contains p arcs whose starting points are $i + 1, \dots, i + p$, whose ending points are $j + 1, \dots, j + p$, and that form the endhered pattern π . Figure 3b shows an example of a matching containing pattern shown on Figure 3a.

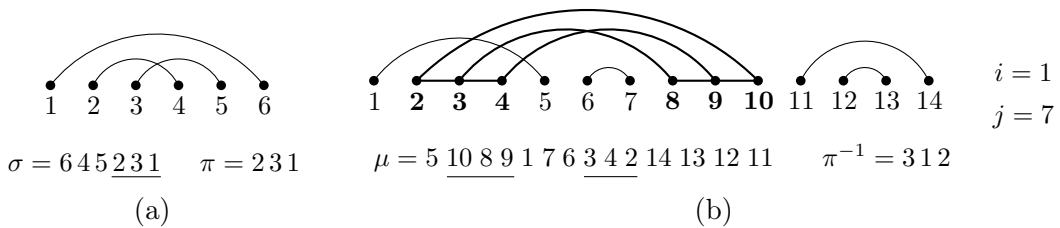


Figure 3: Endhered pattern 231 (a) and an example (b) of its occurrence.

The endhered patterns have been rarely studied in the literature, although they may shed light on the formation of collapsed RNA structures from Rødland’s paper about pseudoknots [41]. The only work in this direction that we are aware of is the paper of Baril [4] who examined one of two endhered patterns of size 2 in his study on irreducible involutions and permutations.

2.2 Endhered twists and engendered symmetries

Given an endhered pattern π , let us denote by $a_{n,k}(\pi)$ the number of matchings of size n with k occurrences of π . If, additionally, τ is another endhered pattern, then we designate by $a_{n,k,m}(\pi, \tau)$ the number of matchings of size n with k and m occurrences of patterns π and τ , respectively. Certain patterns, for instance $\pi = \overbrace{\curvearrowright}$ and $\tau = \overleftarrow{\curvearrowright}$, have the same distribution, meaning that $a_{n,k}(\pi) = a_{n,k}(\tau)$. The goal of this Subsection is to establish such equidistributed classes of endhered patterns with the help of bijections, i.e., without direct enumeration. To this end, we apply matching transformations that we call *endhered twists*.

Definition 2.2. The *left endhered twist* (resp. *right endhered twist*) is a transformation that takes a matching μ and produces a matching $\text{letw}(\mu)$ (resp. $\text{retw}(\mu)$) such that all runs of consecutive starting (resp. ending) points are reversed. Figure 4 shows an example of the right twist.

Observation 2.3. *Endhered twists of endhered patterns correspond to classical symmetries on permutations. Thus, the right endhered twist applied to an endhered pattern $\pi = \pi_1 \dots \pi_p$*

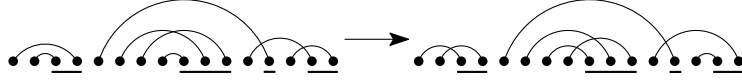


Figure 4: An example of the right endhered twist, runs of right points are underlined.

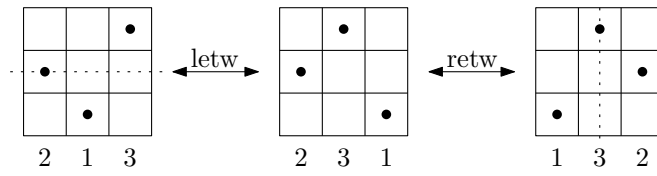


Figure 5: Geometrical meaning of endhered twists.

is its reverse: $\text{retw}(\pi) = \pi_p \dots \pi_1$. At the same time, the left endhered twist is the complement: $\text{letw}(\pi) = (p+1-\pi_1) \dots (p+1-\pi_p)$. For example,

$$\begin{aligned} \text{retw}(321) &= 123, & \text{retw}(231) &= 132, & \text{retw}(312) &= 213, \\ \text{letw}(321) &= 123, & \text{letw}(132) &= 312, & \text{letw}(213) &= 231. \end{aligned}$$


Geometrically, if we represent permutations as square tables, the endhered twists are axial symmetries (see Figure 5).

Lemma 2.4. *Two endhered patterns π and τ of the same size have the same joint distribution if they are identical under the left or right endhered twists. In other words, if $\pi = \text{letw}(\tau)$ or $\pi = \text{retw}(\tau)$, then*




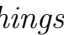
$$a_{n,k,m}(\pi, \tau) = a_{n,k,m}(\tau, \pi) \quad (1)$$

for any integers n , k and m . In particular, $a_{n,k}(\pi) = a_{n,k}(\tau)$.

Proof. Suppose that π is obtained from τ by applying the left endhered twist. Let μ be a matching containing k occurrences of π and m occurrences of τ . Applying the left twist to the whole μ , we transform every occurrence of π to τ and vice versa. No other occurrence of π or τ are created. This implies relation (1). The cases of the right twist is obtained *mutatis mutandis*. \square

Remark 2.5. Given a positive integer n and an endhered pattern π , the value $a_{n,0}(\pi)$ is the number of matchings of size n avoiding π . We say that two endhered patterns π and τ are *Wilf equivalent* if $a_{n,0}(\pi) = a_{n,0}(\tau)$ for every n . For example, Wilf equivalent patterns of size 2 and 3 are shown in Figure 6. Matchings avoiding the pattern  correspond to involutions with no fixed points and no successions considered by Baril [4].

For the endhered patterns of size 2 and 3, as we will see in this paper, their Wilf equivalence implies that the corresponding patterns are equidistributed. In other words, if $a_{n,0}(\pi) = a_{n,0}(\tau)$ for every n , then $a_{n,k}(\pi) = a_{n,k}(\tau)$ for all n and k . In the general case, this question is not trivial.

Corollary 2.6. *The joint distribution of the endhered patterns $\pi = 1 \dots p$ and $\tau = p \dots 1$ is symmetric. In particular, the number of size- n matchings containing k  and m  equals the number of size- n matchings containing m  and k .*

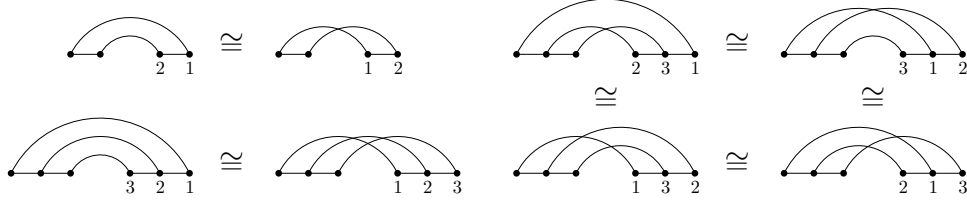


Figure 6: Wilf equivalence classes of endhered patterns of size 2 and 3.

Remark 2.7. Corollary 2.6 is consistent with the result of Chen, Deng, Du, Stanley, and Yan [10] who showed that, in the case of classical patterns, the joint distribution of crossings and nestings is symmetric.

2.3 Enumeration of endhered patterns of size 2

The goal of this Subsection is to study the distributions of endhered patterns of size 2. It follows from Lemma 2.4 that patterns $\overleftarrow{\curvearrowright}$ and $\overrightarrow{\curvearrowright}$ have the same distribution, so we need only to enumerate one of them. For simplicity, we will write $a_{n,k}$ instead of $a_{n,k}(21)$ throughout the rest of the paper. Note that some of results presented here (Lemma 2.9, Corollary 2.10) were proved by Baril [4] using the language of permutations and involutions. For the sake of completeness, we provide new versions of these proofs using the language of matchings.

Theorem 2.8. For $n \geq 1$, the number of size- n matchings containing exactly k occurrences of pattern $\overleftarrow{\curvearrowright}$ (resp. $\overrightarrow{\curvearrowright}$) satisfies

$$\begin{cases} a_{n+1,k} &= a_{n,k-1} + 2(n-k) \cdot a_{n,k} + 2(k+1) \cdot a_{n,k+1} \\ a_{1,0} &= 1 \\ a_{1,k} &= 0, \quad k \neq 0. \end{cases} \quad (2)$$

Formally, we allow negative values of k . It follows from (2) that $a_{n,k} = 0$ whenever $k < 0$.

Proof. The boundary values corresponding to $n = 1$ are easily obtained. To construct a matching of size $n + 1$ having k occurrence of $\overleftarrow{\curvearrowright}$, we take a non-empty matching μ of size n and add a new arc E together with its two ending points (left and right). We always put the left point of E before the first point of μ . For the right point of E , there are three different cases, since the new arc either creates a new pattern, destroys an existing pattern, or leaves the number of occurrences of the pattern unchanged.

1. The right point of E is placed just after the right point of the arc $(1, i)$ of μ (see Figure 7). In this case, the number of occurrences increases by one. Hence, the matching μ must possess $k - 1$ patterns $\overleftarrow{\curvearrowright}$, so that the resulting matching has k occurrences of this pattern. This gives us the first term in relation (2).
2. The right end of E is put between the left (resp. right) ends of any two arcs that make up the pattern $\overleftarrow{\curvearrowright}$. In this case, one of the existing patterns is destroyed, and the number of occurrences decreases by one (see Figure 8). Therefore, to have

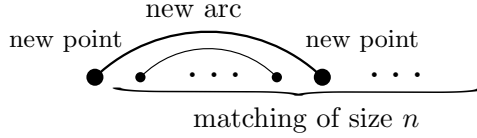


Figure 7: Adding a new arc creates a new pattern.

k occurrences in the resulting matching, the initial matching μ must contain $k + 1$ occurrences of \curvearrowright . Since there are $2(k + 1)$ destructive positions, this gives us the last term in relation (2).

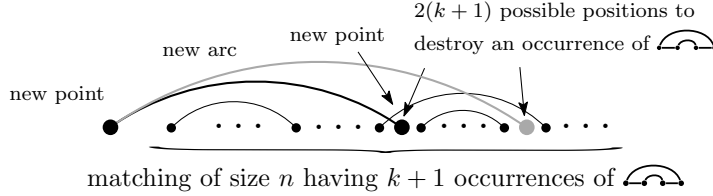


Figure 8: Adding a new arc destroys an existing pattern.

3. The right point of E does not create, nor remove any pattern \curvearrowright . In this case, the matching μ must possess k occurrences of this pattern. There are $2n + 1$ possible positions for the right end of E , but $2k + 1$ of them are forbidden. Indeed, there are k already existing occurrence of \curvearrowright that we do not want to destroy, and one more position is banned because we are not allowed to create any new \curvearrowright in this case (see Figure 9).

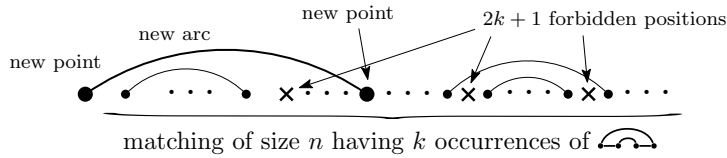


Figure 9: Adding a new arc does not change the number of pattern occurrences.

The distribution of the pattern \curvearrowright respects the same recurrence. To establish the latter fact, it is sufficient to twist the corresponding ends using Lemma 2.4. \square

Table 1 shows first values of $a_{n,k}$. With this table and those that follow in the text, we can see the big picture of the initial terms, intuitively understand how they grow, formulate hypotheses about the global (asymptotic) behaviour of these numerical sequences. We present our results in this direction.

Lemma 2.9. *For any $n > k \geq 0$,*

$$a_{n,k} = \binom{n-1}{k} a_{n-k,0}. \quad (3)$$

$\begin{array}{c} n \\ \backslash \\ k \end{array}$	1	2	3	4	5	6	7	8	9	OEIS
0	1	2	10	68	604	6584	85048	1269680	21505552	A165968
1	0	1	4	30	272	3020	39504	595336	10157440	A179540
2	0	0	1	6	60	680	9060	138264	2381344	
3	0	0	0	1	8	100	1360	21140	368704	
4	0	0	0	0	1	10	150	2380	42280	
5	0	0	0	0	0	1	12	210	3808	
6	0	0	0	0	0	0	1	14	280	
7	0	0	0	0	0	0	0	1	16	
8	0	0	0	0	0	0	0	0	1	

Table 1: Distribution of pattern $\overbrace{\curvearrowright}^k$ (21).

Proof. Each matching of size n with k occurrences of $\overbrace{\curvearrowright}^k$ can be uniquely obtained by the following procedure. Take a matching of size $n - k$ avoiding $\overbrace{\curvearrowright}^k$ pattern. Add k new arcs expanding some of the original arcs into double, triple, \dots , and $(k + 1)$ -uple arcs. Since each t -uple arc contains $t - 1$ occurrences of $\overbrace{\curvearrowright}^k$, we have k occurrences of $\overbrace{\curvearrowright}^k$ in total (see Figure 10). Therefore, relation (3) comes from the classical stars-and-bars argument [14]: the binomial coefficient $\binom{n-1}{k}$ is the number of ways to add k new arcs by expanding some of the $n - k$ original arcs.

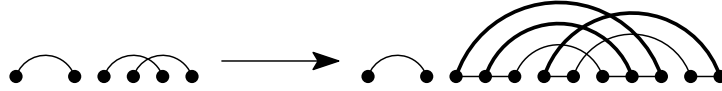


Figure 10: Arc expansion. Adding 3 new arcs, we create 3 new occurrences of $\overbrace{\curvearrowright}^k$.

□

Corollary 2.10. *The sequence $(a_{n,0})$ satisfies the recurrence relation*

$$a_{n+1,0} = 2na_{n,0} + 2(n-1)a_{n-1,0} \quad (4)$$

with the initial conditions

$$a_{1,0} = 1 \quad \text{and} \quad a_{2,0} = 2.$$

In particular, this sequence corresponds to a shift of [A165968](#) in OEIS [35].

Proof. Recurrence (4) follows immediately from relations (2) and (3). □

Remark 2.11. There is a constructive combinatorial explanation of relation (4) (see also the work of Baril [4]). Indeed, if we remove the first arc in a matching of size $n + 1$ avoiding $\overbrace{\curvearrowright}^k$, then we obtain a matching of size n that either avoids $\overbrace{\curvearrowright}^k$ or contains exactly one pattern $\overbrace{\curvearrowright}^k$. This observation shows that any matching of size $n + 1$ avoiding $\overbrace{\curvearrowright}^k$ can be uniquely constructed in the following way. Either we take a matching μ of

size n avoiding $\overbrace{\curvearrowright}$ and add a new arc starting at the left of the already existing $2n$ points of μ and ending at some of $2n$ possible new positions (it cannot arrive just after the right point of the first arc of μ). Or we take a matching of size $n - 1$ avoiding $\overbrace{\curvearrowright}$, double one of its $n - 1$ arcs, and then destroy the just created pattern $\overbrace{\curvearrowright}$ by a new arc drawn in one of two possible ways (see Figure 11).

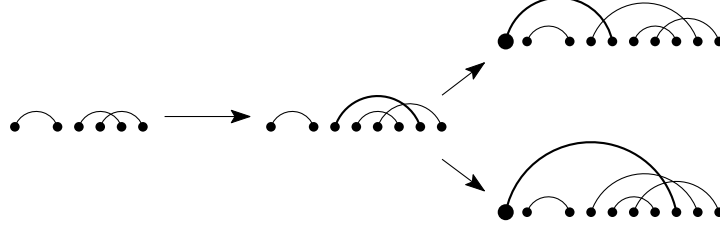


Figure 11: Two possible ways for destroying the newly created pattern $\overbrace{\curvearrowright}$.

Corollary 2.12. *The sequence $(a_{n,1})$ satisfies the recurrence relation*

$$a_{n+1,1} = 2n(a_{n,1} + a_{n-1,1}) \quad (5)$$

with the initial conditions

$$a_{1,1} = 0 \quad \text{and} \quad a_{2,1} = 1.$$

In particular, this sequence corresponds to a shift of [A179540](#) in OEIS [35].

Proof. Apply relations (2) and (3):

$$a_{n+1,1} = na_{n,0} = n(2(n-1)a_{n-1,0} + 2a_{n-1,1}) = 2n(a_{n,1} + a_{n-1,1}).$$

□

Lemma 2.13. *For any $n \geq 0$,*

$$a_{n+1,0} = \sum_{k=0}^n (-1)^{n-k} \binom{n}{k} (2k+1)!!, \quad (6)$$

where $m!!$ denotes the double factorial of m ,

$$m!! = \begin{cases} m(m-2)(m-4) \cdots 4 \cdot 2 & \text{if } m \text{ is even,} \\ m(m-2)(m-4) \cdots 3 \cdot 1 & \text{if } m \text{ is odd.} \end{cases}$$

Proof. Being rewritten as

$$a_{n+1,0} = \sum_{k=0}^n (-1)^k \binom{n}{k} (2(n-k)+1)!!,$$

relation (6) can be proved with the help of the inclusion-exclusion principle. Indeed, there are $(2n+1)!!$ matchings of size $n+1$. In $(2n-1)!!$ among them, i -th arc belongs to the pattern $\overbrace{\quad}$ as an outer arc, $i=1, \dots, n$. Hence, to get the number of pattern avoiding matchings, we need to deduce $\binom{n}{1}(2n-1)!!$ from $(2n+1)!!$. However, matchings with two occurrences of $\overbrace{\quad}$ are deduced twice. There are $(2n-3)!!$ matchings with i -th and j -th arcs belonging to the pattern $\overbrace{\quad}$ as outer arcs, and $\binom{n}{2}$ possible pairs (i, j) , so we add $\binom{n}{2}(2n-3)!!$. The same way, we treat matchings with three occurrences of $\overbrace{\quad}$, with four occurrences and so on. □

For any $n, k \geq 0$, let us denote

$$b_{n,k} := a_{n+1,k}.$$

Define also

$$B(z, u) := \sum_{n=0}^{\infty} \sum_{k=0}^n b_{n,k} \frac{z^n}{n!} u^k.$$

Lemma 2.14. *The exponential generating function $B(z, u)$ satisfies*

$$B(z, u) = \frac{e^{z(u-1)}}{\sqrt{(1-2z)^3}}. \quad (7)$$

In particular, the exponential generating function of the shifted k -th row of Table 1 is

$$[u^k]B(z, u) = \frac{z^k}{k!} \cdot \frac{e^{-z}}{\sqrt{(1-2z)^3}}.$$

Proof. Recall that

$$\sum_{n=0}^{\infty} (2n+1)!! \frac{z^n}{n!} = \frac{1}{\sqrt{(1-2z)^3}}.$$

Therefore, according to Lemma 2.13, we have

$$\begin{aligned} B(z, 0) &= \sum_{n=0}^{\infty} b_{n,0} \frac{z^n}{n!} = \sum_{n=0}^{\infty} \frac{z^n}{n!} \sum_{k=0}^n (-1)^{n-k} \binom{n}{k} (2k+1)!! \\ &= \left(\sum_{n=0}^{\infty} \frac{(-z)^n}{n!} \right) \left(\sum_{n=0}^{\infty} (2n+1)!! \frac{z^n}{n!} \right) = \frac{e^{-z}}{\sqrt{(1-2z)^3}}. \end{aligned}$$

Hence, due to Lemma 2.9,

$$\begin{aligned} B(z, u) &= \sum_{n=0}^{\infty} \frac{z^n}{n!} \sum_{k=0}^n \binom{n}{k} b_{n-k,0} u^k \\ &= \left(\sum_{n=0}^{\infty} b_{n,0} \frac{z^n}{n!} \right) \left(\sum_{n=0}^{\infty} \frac{(zu)^n}{n!} \right) = \frac{e^{-z}}{\sqrt{(1-2z)^3}} \cdot e^{zu}. \end{aligned}$$

□

Remark 2.15. Combinatorially, relation (7) represents another facet of the inclusion-exclusion principle. This can be proved using the symbolic method, see [15, Section III.7.4]. Indeed, the exponential generating function $e^{zv}/\sqrt{(1-2z)^3}$ corresponds to the labeled product of the combinatorial class of urns \mathcal{U} and the derivative of the combinatorial class of matchings \mathcal{M} . Urns can be interpreted as distinguished arcs (marked by the variable v) that are inserted into a matching to form double arcs, that is, patterns $\overbrace{\curvearrowright}$. Passing to the derivative \mathcal{M}' corresponds to adding a supplementary arc to all matchings. This operation guarantees that every matching possesses at least one arc (and hence, there is something to double). Finally, the variable substitution $v = u - 1$ has the meaning of the inclusion-exclusion itself.

2.4 Asymptotics for patterns of size 2

Here, we provide the asymptotics of the distribution of size-2 patterns. These formulas allow us to understand in a concise way how the pattern distribution behaves when the size of matchings increases, without the need for precise calculations.

Theorem 2.16. *The asymptotic behavior of the numbers of matchings with k occurrences of pattern $\overbrace{\curvearrowright}$ (resp. $\overbrace{\curvearrowleft}$), as $n \rightarrow \infty$, is*

$$a_{n,k} \sim \frac{1}{2^k k!} \left(\frac{2}{e}\right)^{n+1/2} n^n.$$

Proof. Let us start with establishing the asymptotics of $b_{n,0}$. Due to Lemma 2.13, we have

$$\begin{aligned} b_{n,0} &= \sum_{k=0}^n (-1)^{n-k} \binom{n}{k} (2k+1)!! \\ &= (2n+1)!! \sum_{k=0}^n \frac{(-1)^k}{k!} \cdot \frac{n(n-1)\dots(n-k+1)}{(2n+1)(2n-1)\dots(2n-2k+3)} \\ &\sim (2n+1)!! \sum_{k=0}^n \frac{(-1)^k}{2^k k!} \sim \frac{(2n)(2n)!}{2^n n! e^{1/2}} \sim \frac{2^{n+3/2}}{e^{n+1/2}} \cdot n^{n+1}, \end{aligned}$$

where the last passage is done with the help of Stirling's formula. Now, the general case can be done by Lemma 2.9:

$$b_{n,k} = \binom{n}{k} b_{n-k,0} \sim \frac{n^k}{k!} \cdot \frac{2^{n-k+3/2}}{e^{n-k+1/2}} \cdot (n-k)^{n-k+1} \sim \frac{2^{n-k+3/2}}{k!} \cdot \frac{n^{n+1}}{e^{n+1/2}}.$$

Passing to $a_{n,k} = b_{n-1,k}$, we get the final result. \square

Remark 2.17. There is an alternative proof that relies on the singularity analysis of the exact form of the series $B(z, u)$ given by Lemma 2.14. For details of this technique, see [15, Theorem VI.4].

Corollary 2.18. *The limit distribution of pattern $\overbrace{\curvearrowright}$ (resp. $\overbrace{\curvearrowleft}$) in a uniform random matching of size n follows asymptotically a Poisson law with parameter $1/2$.*

Proof. This follows from Theorem 2.16 and Stirling's formula:

$$\frac{a_{n,k}}{(2n-1)!!} \sim \frac{e^{-1/2}}{2^k k!}.$$

□

Corollary 2.19. *The ratio of numbers from the k -th row to the numbers of the $(k+1)$ -th row of Table 1 tends to $2(k+1)$. In other words, for any $k \in \mathbb{N}$,*


$$\frac{a_{n,k}}{a_{n,k+1}} \sim 2(k+1).$$

Proof. This is an immediate consequence of Corollary 2.18. □


2.5 Patterns of size 3

In this Subsection, we examine endhered patterns of size 3. There are six of them and, as it follows from Lemma 2.4, they are divided into two equivalence classes with respect to distributions of these patterns in matchings (see Fig. 6). Thus, our goal is to study these two different distributions represented, for instance, by patterns 321 and 132. For simplicity, we will use the following notations:

$$c_{n,k} := a_{n,k}(321) = a_{n,k}(123)$$


to denote the number of matchings of size n containing k patterns 321 () , and



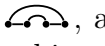

$$d_{n,k} := a_{n,k}(132) = a_{n,k}(213) = a_{n,k}(231) = a_{n,k}(312)$$




to designate the number of matchings of size n containing k patterns 132 ().

Theorem 2.20. *For any $n > 0$ and $k \geq 0$,*

$$\begin{aligned} c_{n,0} &= \sum_{s=0}^{\lfloor n/2 \rfloor} \binom{n-s}{s} a_{n-s,0}, \\ c_{n,k} &= \sum_{s=1}^{\lfloor (n-k)/2 \rfloor} \binom{k+s-1}{k} \binom{n-k-s}{s} a_{n-k-s,0} \quad \text{if } k > 0, \end{aligned} \tag{8}$$

where $a_{n,0}$ is the number of matchings of size n avoiding pattern 21 ().

Proof. Each matchings avoiding pattern  can be uniquely obtained from a matching avoiding pattern  by doubling some of its arcs (possibly, none). In particular, to get a matching of size n by doubling exactly s arcs, where $0 \leq s \leq \lfloor n/2 \rfloor$, we take a matching of size $n-s$ avoiding , and then choose s of its arcs to transform them into s occurrences of . Since this can be done in $\binom{n-s}{s}$ ways, we obtain relation (8) for $k = 0$.

Similarly, any matching of size n with $k > 0$ occurrences of pattern  can be uniquely obtained by the following procedure. First, we fix $0 < s \leq \lfloor (n-k)/2 \rfloor$ and take a matching of size $n-k-s$ avoiding . Second, we choose s arcs in this matching and double them, there are $\binom{n-k-s}{s}$ ways to do it. Finally, we stack additional k arcs onto just created s occurrences of , which can be done in $\binom{k+s-1}{k}$ ways (apply the stars-and-bars argument [14]). □

$\begin{array}{c} n \\ \backslash \\ k \end{array}$	1	2	3	4	5	6	7	8	9
0	1	3	14	100	906	10022	130864	1969884	33583700
1	0	0	1	4	34	332	3866	52400	811248
2	0	0	0	1	4	36	362	4304	59256
3	0	0	0	0	1	4	38	392	4752
4	0	0	0	0	0	1	4	40	422
5	0	0	0	0	0	0	1	4	42
6	0	0	0	0	0	0	0	1	4
7	0	0	0	0	0	0	0	0	1



Table 2: Distribution of pattern  (321).

Table 2 shows the first values of $c_{n,k}$.

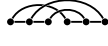
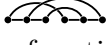
In order to study the distribution of the occurrences of pattern , we apply Goulden-Jackson *cluster method* [19, 20] and the inclusion-exclusion principle in the framework of the symbolic method [15, Section III.7.4]. Let us trace the main steps.

We start by considering all possible matchings; the corresponding generating function is given by

$$F(z) = \sum_{n=0}^{\infty} (2n-1)!! z^n = 1 + z + 3z^2 + 15z^3 + 105z^4 + \dots$$

In every matching, we distinguish certain arcs, say, by coloring them violet. Algebraically, this is done by passing to the generating function


$$G(z, v) = \sum_{n=0}^{\infty} \sum_{k=0}^n g_{n,k} z^n v^k = F(z + zv),$$

where $g_{n,k}$ is the number of matchings of size n with k violet arcs (marked by the variable v). Next, we replace violet arcs by “thick arcs” constructed from 3 arcs forming pattern  (in general case, we should replace them by *clusters* consisting of intersected copies of a given pattern; pattern , however, admits no possible self-intersection). This corresponds to the generating function

$$H(z, v) = G(z, z^2v).$$

And finally, according to the inclusion-exclusion principle, we obtain the generating function

$$D(z, u) = \sum_{n=0}^{\infty} \sum_{k=0}^n d_{n,k} z^n u^k = H(z, u-1)$$

whose coefficients $d_{n,k}$ enumerate the matchings of size n having k occurrences of the pattern  (see Figure 12 illustrating the process).

Thus, the following result takes place.

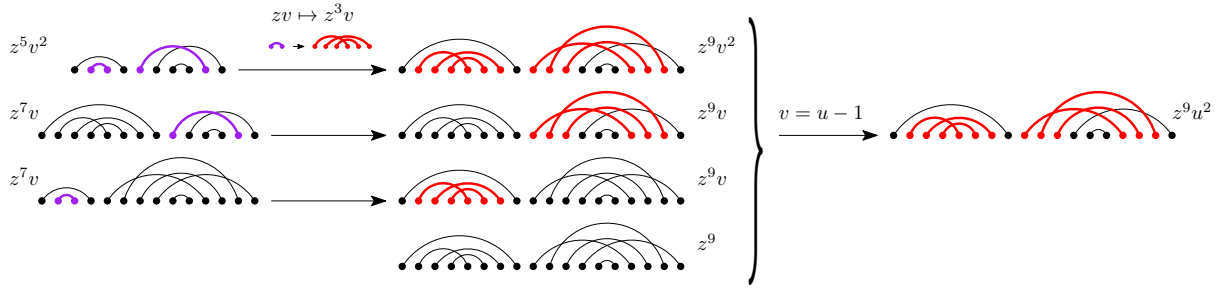



Figure 12: Cluster method and the inclusion-exclusion principle.


Theorem 2.21. *The bivariate generating function*

$$D(z, u) = \sum_{n=0}^{\infty} \sum_{k=0}^n d_{n,k} z^n u^k,$$

where $d_{n,k}$ is the number of matchings of size n containing k patterns  can be expressed as

$$D(z, u) = F(z + (u - 1)z^3) = \sum_{n=0}^{\infty} (2n - 1)!! (z + (u - 1)z^3)^n.$$

$\begin{array}{c} n \\ k \end{array}$	1	2	3	4	5	6	7	8	9
0	1	3	14	99	900	9978	130455	1965285	33522915
1	0	0	1	6	45	414	4635	61110	927090
2	0	0	0	0	0	3	45	630	9405
3	0	0	0	0	0	0	0	0	15

Table 3: Distribution of pattern  (132).




The limit distributions of patterns  and  are different. In the case of , for a fixed $k \geq 0$, as $n \rightarrow \infty$, we have

Table 3 shows the first values of $d_{n,k}$.

$$c_{n,k} \sim \sqrt{2} \cdot C_k \cdot \frac{(2n)^{n-k}}{e^n},$$

where

$$C_0 = 1 \quad \text{and} \quad C_k = \sum_{s=1}^k \binom{k-1}{s-1} \frac{1}{2^s s!} \quad \text{if } k > 0.$$

On the other hand, for  takes place

$$d_{n,k} \sim \frac{2^{n-2k+1/2}}{k! e^n} \cdot n^{n-k},$$

as $n \rightarrow \infty$. In other words, for large values of n , we have

$$\frac{c_{n,k}}{(2n-1)!!} \sim \frac{C_k}{2^k} \cdot \frac{1}{n^k} \quad \text{and} \quad \frac{d_{n,k}}{(2n-1)!!} \sim \frac{1}{2^{2k}k!} \cdot \frac{1}{n^k},$$

meaning that a large uniform random matching avoids both patterns with high probability. The proof of this asymptotic behavior requires other technical means and is beyond the scope of this paper. We will present this proof as part of the analysis of the asymptotic behavior of any (more complex) endhered pattern in our next work.

3 Occurrences of endhered patterns in RNA secondary structures and shapes

In this section, we shift our attention to the native (real-world) data and discuss how endhered patterns are distributed in various representations of the secondary RNA structures obtained from Protein Data Bank [6] (PDB, <https://www.rcsb.org>) using our Python scripts (https://gitlab.com/celiabiane/endhered_pattern), which depends on GEMMI [47] to parse mmCIF files from PDB (<https://github.com/project-gemmi/gemmi>), and uses two methods to detect base pairs in RNA molecules: a closed-source software x3DNA-DSSR [31,32] and an open-source software FR3D-python [36,42] (<https://github.com/BGSU-RNA/fr3d-python>).

The software x3DNA-DSSR directly produces extended dot-bracket notations. To derive these notations x3DNA-DSSR uses only *canonical pairs*: A-U, C-G, wobble G-U, and A-T (in RNA-DNA hybrids) with *cis*. Watson-Crick/Watson-Crick interactions and without forming parallel mini-duplexes. FR3D-python gives a list of base pairs. We parse its results, filter canonical base pairs *à la* x3DNA-DSSR, and produce extended dot-bracket notations using a simple First-Come-First-Served method.

There are several hundred known, existing in nature, modifications of nucleotides. In the data they are denoted by special one-, two-, or three-letter codes, different from the classical 4 letters UACG. Modified nucleotides are mapped to short 1-letter nucleotide names. For example A23 is mapped A, EQ0 to G, and CCC to C. Nucleic Acid KnowledgeBase [5, 29] (<https://www.nakb.org/modifiednt.html>) contains details for these mappings. Following x3DNA-DSSR¹ we adapt one exception to these rules: pseudouridine (PSU) is mapped to P and not to U. This adaptation allows us to better compare the x3DNA-DSSR and FR3D-python results.

From the RNA secondary structures obtained with x3DNA-DSSR and FR3D-python, we collapsed unpaired nucleotides in order to keep only paired ones. When keeping only RNA structures composed of one chain, this leads to 1501 RNA structures, in 933 (resp. 929) of them x3DNA-DSSR (resp. FR3D-python) has found at least one canonical base pair. The data has been accessed on August 29, 2024.

The structures in the extended dot-bracket notation are converted to matchings using an algorithm based on stacks. The algorithm works as follows: the word composed of parentheses is read from left to right, when an opening character is met its index is stacked in a stack corresponding to the nature of the character. When a closing character is met,

¹See <http://forum.x3dna.org/rna-structures/modified-nucleotides-incorrect>

a pair corresponding to the last index in the corresponding stack and the current index is added to the matching, and the index of the opening parenthesis is unstacked. Figure 13 shows an example of this conversion.

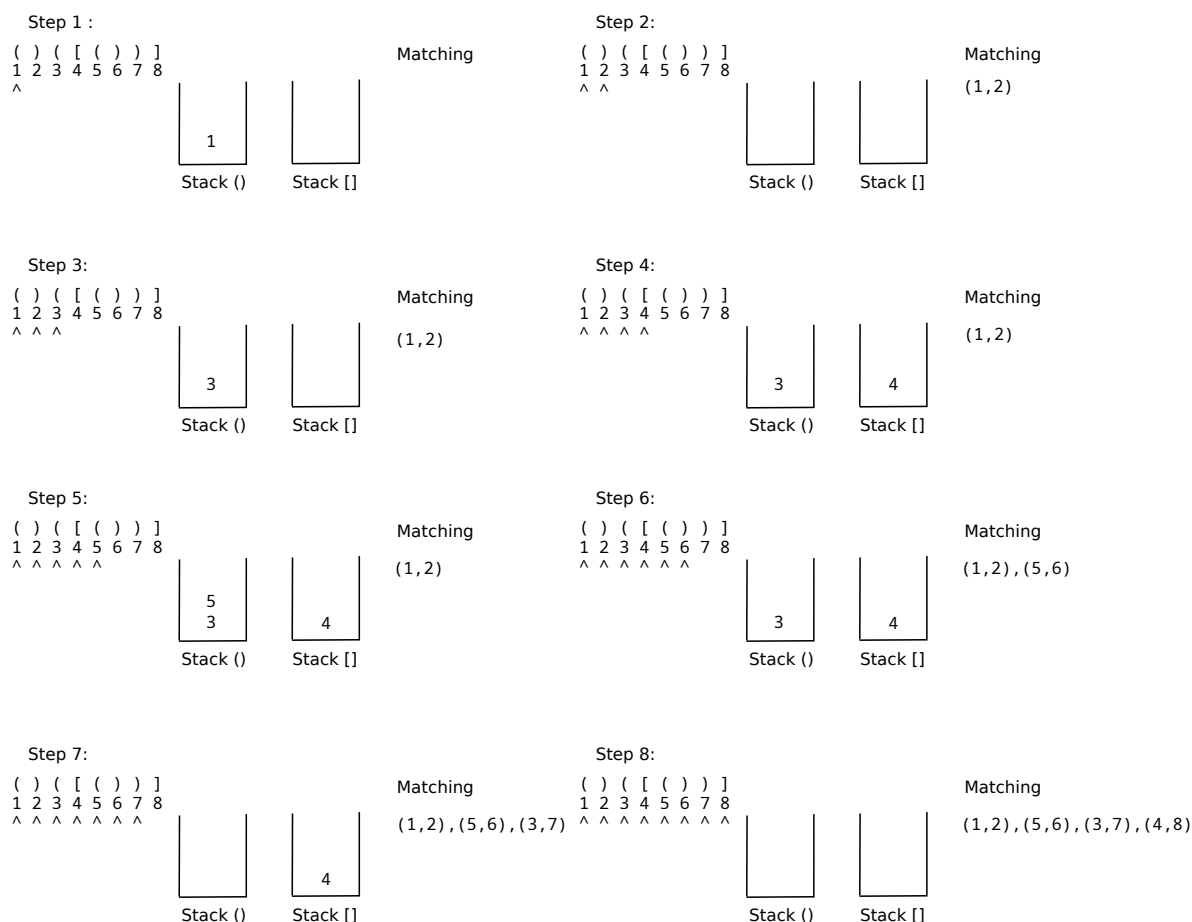




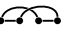
Figure 13: Example of conversion from dot-bracket notation to a matching.

Among the 933 (929 in case of FR3D-python) paired and non empty secondary structures, 926 (921 in case of FR3D-python) contains the pattern $\overbrace{()}$. In case of x3DNA-DSSR, the seven remaining structures, with PDB identifiers 3IZY, 5UZ9, 6B45, 6B46, 6B47, 6B48, 6VQX, contain a single base pairing. In case of FR3D-python, there are 12 of them: 1FC8, 1R4H, 3IZY, 4R8I, 4Z7K, 5UZ9, 6B45, 6B46, 6B47, 6B48, 6BY5, 6VQX. A detailed comparison of the results of the two programs would be very enriching, but is beyond the scope of this article.

None of the considered 933 structures are encoded with more than 4 types of parentheses. In case of x3DNA-DSSR, only 11 structures (4DS6, 4E8V, 8OLS, 8OLV, 8OLW, 8OM0, 8RUI, 8RUJ, 8RUL, 8RUM, 8RUN) are written using 4 types of parentheses, 27 employ 3 types of parentheses (5KMZ, 5TPY, 6AGB, 6AHR, 6UES, 6WLQ, 6WLR, 6WLS, 7EZ0, 7K16, 7U4A, 7UO5, 7UVT, 7XD3, 7XD4, 7XD7, 7XSK, 7XSL, 7XSM, 7XSN, 8HD6, 8HD7, 8I7N, 8T29, 8T2A, 8T2B, 8T2O), and 193 structures are encoded with 2 types of parentheses. In case of FR3D-python, the corresponding numbers are 1, 39 and 184. We need to be careful when comparing these numbers, x3DNA-DSSR uses elimination-gain heuristics [43] to produce extended dot-bracket notations, while we adapt

First-Come-First-Served method.

The distribution of patterns  and  in RNA secondary structure is shown in Figure 15. We can observe that the majority of RNA secondary structures are of small size (between 0 and 100 nucleotides), and there seems to be a linear relation between the size of the matching and the number of occurrence of the pattern. This can be explained by the large number of hairpin patterns in RNA secondary structures.

Both x3DNA-DSSR and FR3D-python detect 2 RNA structures containing the  pattern: 4M4O, 5U3G. One of them, 4M4O, corresponds to the minE aptamer involved in a complex with a lysozyme. Both methods, based on x3DNA-DSSR and FR3D-python, give the following secondary structure:

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

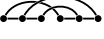
Another one, 5U3G, is the *Dickeya dadantii* ykkC riboswitch, which has the following secondary structure:




x3DNA-DSSR: ((((((.....((((([..])..)))))))).))))(((((.....))))).([...((.....))..)..
 FR3D-python: .((((.....((((([..])..)))))))).((((.....))))).([...((.....))..)..

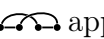
Diving into the data, we see that FR3D-python does not detect GTP-C pair on position (1,40) in this case.

The results of two methods differ in the case of molecule 7K16 (Tamana Bat Virus xrRNA1):

x3DNA-DSSR: {{...(((((((.....))))))(((((.....))))))...}}
 FR3D-python: .(...[[(((((((.....))))))(((((.....))))))...}}

In this case x3DNA-DSSR detects, in addition to FR3D-python, a 5GP-C base pair on position (1,42). Removing this pair, we create a  pattern. See Table 4 for more details.

Pattern  is contained only in 4M4O, while pattern  belongs only to 5U3G. Pattern  is contained in 918 (908 in case of FR3D-python) structures.

It is surprising that pattern  appears so rarely in RNA structures while pseudoknots are thought to have important biological functions. We observe that occurrence of this pattern are “hidden” by the high frequency of nested bonds (see Figure 14). To neutralize this effect, we pass to RNA shapes.

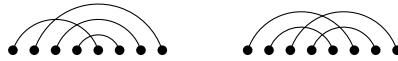




Figure 14: Two examples of matchings with nesting bonds (occurrences of ) but without occurrences of .

The concept of RNA shapes have been introduced by Giegerich, Voss, and Rehmsmeier in 2004 [17]. In these shapes, no unpaired regions are included and nested bonds are combined. For instance, the secondary structure

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

has the following RNA shape:

((()))


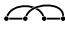




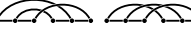

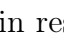
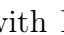
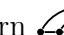

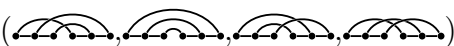



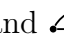
Pattern	RNA secondary structure		Extended RNA shape	
	x3DNA-DSSR	FR3D-python	x3DNA-DSSR	FR3D-python
	926	921	0	0
	2 (4M4O, 5U3G)	3 (4M4O, 5U3G, 8SH5)	59 (1A60, 1E95, 1KAJ, 1KPD, 1KPY, 1KPZ, 1L2X, 1L3D, 1RNK, 1YG3, 1YG4, 1YMO, 2A43, 2AP0, 2AP5, 2G1W, 2K95, 2K96, 2LC8, 2M58, 2M8K, 2RP0, 2RP1, 2TPK, 2XDB, 3IYQ, 3IYR, 3IZ4, 3U4M, 3U56, 3UMY, 437D, 4M4O, 4PQV, 4QG3, 4QVI, 4R8I , 5KMZ, 5NPM, 5TPY, 5U3G, 6AGB, 6AHR, 6DLQ, 6DLR, 6DLS, 6DLT, 6DNR, 6E1T, 6E1V, 6VUH, 7K16, 7U4A, 7UO5, 8HIO, 8T29, 8T2A, 8T2B, 8T2O)	60 (1A60, 1E95, 1KAJ, 1KPD, 1KPY, 1KPZ, 1L2X, 1L3D, 1RNK, 1YG3, 1YG4, 1YMO, 2A43, 2AP0, 2AP5, 2G1W, 2K95, 2K96, 2LC8, 2M58, 2M8K, 2RP0, 2RP1, 2TPK, 2XDB, 3IYQ, 3IYR, 3IZ4, 3U4M, 3U56, 3UMY, 437D, 4M4O, 4PQV, 4QG3, 4QVI, 5KMZ, 5NPM, 5TPY, 5U3G, 6AGB, 6AHR, 6D3P , 6DLQ, 6DLR, 6DLS, 6DLT, 6DNR, 6E1T, 6E1V, 6VUH, 7K16, 7U4A, 7UO5, 8HIO, 8SH5 , 8T29, 8T2A, 8T2B, 8T2O)
	1 (4M4O)	1 (4M4O)	8 (3U4M, 3U56, 3UMY, 4M4O, 4QG3, 4QVI, 4R8I , 5NPM)	7 (3U4M, 3U56, 3UMY, 4M4O, 4QG3, 4QVI, 5NPM)
	1 (5U3G)	1 (5U3G)	13 (3U4M, 3U56, 3UMY, 4QG3, 4QVI, 5KMZ, 5NPM, 5U3G, 6DLQ, 6DLR, 6DLS, 6DLT, 6DNR)	13 (3U4M, 3U56, 3UMY, 4QG3, 4QVI, 5KMZ, 5NPM, 5U3G, 6DLQ, 6DLR, 6DLS, 6DLT, 6DNR)
	0	1 (7K16)	0	0
	918	908	0	0
	0	0	0	0

Table 4: RNA secondary structures and RNA shapes in which a given pattern occurs at least once. Differences are highlighted in bold.

Originally, the RNA shapes have been defined in words with a single type of parenthesis. They are counted by Motzkin numbers [12, 13, 17] and correspond exactly to non-crossing matchings avoiding the endhered pattern .

We adapt the Giegerich-Voss-Rehmsmeier reduction to matchings with crossings represented by words with different types of parentheses. The result of this adaptation is what Rødland called *collapsed structures* [41]. This is done by keeping only (i, j) pairs in matching such that $(i + 1, j - 1)$ does not belong to the matching and then reindexing the pairs. The number of  patterns in resulting reduced matchings (RNA shapes) is obviously 0. Interestingly, the number of RNA shapes with at least one occurrence of pattern  increases up to 59 (60 with FR3D-python). Among those, 8 (7 with FR3D-python) RNA structures have pattern  and 13 have pattern . Other size 3 patterns () are not detected in RNA shapes. It is expected for patterns , , and , as they contain the pattern  forbidden in RNA shapes.

In this paper, we study only endhered patterns of size 2 and 3. The analysis of more complex patterns in native (real-world) data may be a part of future works. The function `count_pattern` in `count_visualisation.py` from our GitLab repository can be applied to any endhered pattern.

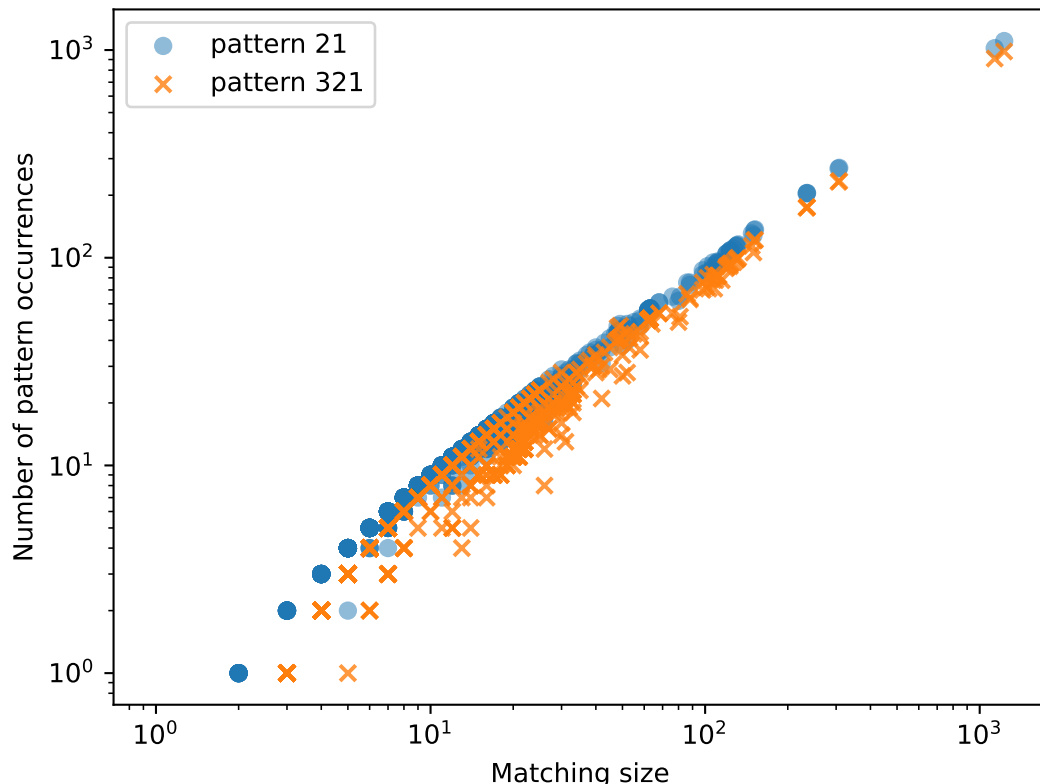


Figure 15: Scatter plot of the number of occurrences of patterns 21 and 321 as a function of the size of the matching. Axes are in logarithmic scale. RNA with no occurrences of patterns 21 and 321 are not displayed. Matchings are obtained using FR3D-python, the results for x3DNA-DSSR are similar and available on git repository https://gitlab.com/celiabiane/endhered_pattern.


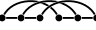
4 Conclusion and discussions

We have examined distributions of endhered patterns of sizes 2 and 3 in matchings from theoretico-combinatorial and data-driven points of view. In matchings, patterns $\overbrace{..}$ and $\overbrace{.}$ have the same distribution. Six endhered patterns of size 3 are divided into 2 equidistributed groups: $\overbrace{..}$, $\overbrace{.}$ and $\overbrace{..}$, $\overbrace{.}$, $\overbrace{..}$, $\overbrace{.}$. Moreover, the joint distribution of patterns of the same size is symmetric if the patterns are equivalent under the twist operation. We have also provided corresponding asymptotic behavior of these distributions.

In this work, we deliberately abstract from the nucleic acid sequences, and model the secondary structures directly by matchings. Our results show that there is a big difference between observed and modelled pattern distributions. This means that non-restricted matchings are too permissive and new models should be developed to get closer to the observed pattern distributions. We wonder if it is possible to describe essential features of RNA secondary structures with pseudoknots, using pattern-based

restrictions. The classical Waterman’s definition of RNA structures and its generalizations (see Subsection 1.1) can be regarded as an example of pattern-based restriction used, among other things, to control the prediction algorithms for secondary structures from nucleic acid sequences. More insights into patterns frequencies in the native secondary RNA structures may guide us towards mixture of two definitions, complex enough to cover different pseudoknot-like structures presented in real data, but at the same time quite simple and neat to prevent the uncontrolled combinatorial explosion.

The non-existence of certain short sequences in genomic and protein sequences is a well-known fact [23, 25, 28, 37]. Applications include cancer research [1, 34] and forensic science [18]. Less is known about the forbidden secondary structures, although some interesting works about theoretical (im)possibility of inverse RNA folding have been published recently [48, 49]. One of the following research directions would be to determine what influence the distribution of endhered patterns in the native RNA. Some configurations are probably forbidden due to physicochemical constraints on the bending of RNA (something like Waterman-Ponty restrictions, but for the case that include pseudoknots), others may not be present because of biological reasons. Are there some evolutionary mechanisms that divert the distribution of patterns from the theoretically observed in the equi-probabilistic model of matchings? How pattern distributions in secondary structures are related to RNA dynamics and function?

For any new combinatorial characterization of RNA structures, we need to develop a method for estimating their affinity with structures observed in native molecules. One such method could also be pattern-based: compare the distribution of patterns in native RNA with theoretically calculated distributions over matchings avoiding certain patterns. Our exploratory study suggest, for instance, that the patterns  and  never appear in RNA. Moreover, PDB references presented in Table 4 look very interesting, especially 4M4O, 5U3G, and 7K16.

Acknowledgments

We would like to express our immeasurable gratitude to Matteo Cervetti, Yann Ponty for the discussions they had with us during the beginning of the work on the endhered patterns, to Justin Masson for help with python code, to Daniel Pinson for proofreading the article, and to the anonymous reviewers for their valuable comments and suggestions. Authors were supported in part by ANR-22-CE48-0002 funded by l’Agence Nationale de la Recherche and the project ANER ARTICO funded by Bourgogne-Franche-Comté region (France).

References

- [1] Nilufar Ali, Cody Wolf, Swarna Kanchan, Shivakumar R. Veerabhadraiah, Laura Bond, Matthew W. Turner, Cheryl L. Jorcyk, and Greg Hampikian. 9S1R nullomer peptide induces mitochondrial pathology, metabolic suppression, and enhanced immune cell infiltration, in triple-negative breast cancer mouse model. *Biomedicine & Pharmacotherapy*, 170:115997, 2024.

- [2] Jørgen Ellegaard Andersen, Josef Mattes, and Nicolai Reshetikhin. The poisson structure on the moduli space of flat connections and chord diagrams. *Topology*, 35(4):1069–1083, 1996.
- [3] Jörg Arndt. *Matters Computational: Ideas, Algorithms, Source Code*. Springer, Berlin, Heidelberg, 2011.
- [4] Jean-Luc Baril. Avoiding patterns in irreducible permutations. *Discrete Mathematics & Theoretical Computer Science*, 17(3):2158, 2016.
- [5] Helen M. Berman, Catherine L. Lawson, and Bohdan Schneider. Developing community resources for nucleic acid structures. *Life*, 12(4):540, 2022.
- [6] Helen M. Berman, John Westbrook, Zukang Feng, Gary Gilliland, T. N. Bhat, Helge Weissig, Ilya N. Shindyalov, and Philip E. Bourne. The Protein Data Bank. *Nucleic Acids Research*, 28(1):235–242, 2000.
- [7] Jonathan Bloom and Sergi Elizalde. Pattern avoidance in matchings and partitions. *The Electronic Journal of Combinatorics*, 20(2), 2013.
- [8] Rutwig Campoamor-Stursberg and Vassily O. Manturov. Invariant tensor formulas via chord diagrams. *Journal of Mathematical Sciences*, 128(4):3018–3029, 2005.
- [9] Matteo Cervetti and Luca Ferrari. Enumeration of some classes of pattern avoiding matchings, with a glimpse into the matching pattern poset. *Annals of Combinatorics*, 26(4):971–995, 2022.
- [10] William Y.C. Chen, Eva Y.P. Deng, Rosena R.X. Du, Richard P. Stanley, and Catherine H. Yan. Crossings and Nestings of Matchings and Partitions. *Transactions of the American Mathematical Society*, 359(4):1555–1575, 2007.
- [11] William Y.C. Chen, Toufik Mansour, and Sherry H.F. Yan. Matchings avoiding partial patterns. *The Electronic Journal of Combinatorics*, 13(1):R112, 2006.
- [12] Sang Kwan Choi, Chaiho Rim, and Hwajin Um. Narayana Number, Chebyshev Polynomial and Motzkin Path on RNA Abstract Shapes. In Jan de Gier, Cheryl E. Praeger, and Terence Tao, editors, *2017 MATRIX Annals*, MATRIX Book Series, pages 153–166. Springer International Publishing, Cham, 2019.
- [13] Marie-Pierre Delest and Gérard Viennot. Algebraic languages and polyominoes enumeration. *Theoretical Computer Science*, 34(1):169–206, 1984.
- [14] William Feller. *An Introduction to Probability Theory and its Applications. Volume II*. John Wiley & Sons, 1991.
- [15] Philippe Flajolet and Robert Sedgewick. *Analytic Combinatorics*. Cambridge University Press, 2009.
- [16] Hin Hark Gan, Samuela Pasquali, and Tamar Schlick. Exploring the repertoire of RNA secondary motifs using graph theory; implications for RNA design. *Nucleic Acids Research*, 31(11):2926–2943, 2003.

- [17] Robert Giegerich, Björn Voss, and Marc Rehmsmeier. Abstract shapes of RNA. *Nucleic Acids Research*, 32(16):4843–4851, 2004.
- [18] Jayita Goswami, Michael C. Davis, Tim Andersen, Abdelkrim Alileche, and Greg Hampikian. Safeguarding forensic DNA reference samples with nullomer barcodes. *Journal of Forensic and Legal Medicine*, 20(5):513–519, 2013.
- [19] Ian P. Goulden and David M. Jackson. An inversion theorem for cluster decompositions of sequences with distinguished subsequences. *Journal of the London Mathematical Society*, s2-20(3):567–576, 1979.
- [20] Ian P. Goulden and David M. Jackson. *Combinatorial Enumeration*. Wiley, 1983.
- [21] Sam Griffiths-Jones, Alex Bateman, Mhairi Marshall, Ajay Khanna, and Sean R. Eddy. Rfam: an RNA family database. *Nucleic Acids Research*, 31(1):439–441, 2003.
- [22] Sam Griffiths-Jones, Simon Moxon, Mhairi Marshall, Ajay Khanna, Sean R. Eddy, and Alex Bateman. Rfam: annotating non-coding RNAs in complete genomes. *Nucleic Acids Research*, 33(suppl_1):D121–D124, 2004.
- [23] Greg Hampikian and Tim Andersen. Absent sequences: nullomers and primes. In *Biocomputing*. World Scientific, 2007.
- [24] Christian Haslinger and Peter F. Stadler. RNA structures with pseudo-knots: graph-theoretical, combinatorial, and statistical properties. *Bulletin of Mathematical Biology*, 61:437–467, 1999.
- [25] Julia Herold, Stefan Kurtz, and Robert Giegerich. Efficient computation of absent words in genomic sequences. *BMC Bioinformatics*, 9(1), 2008.
- [26] Vít Jelínek. Dyck paths and pattern-avoiding matchings. *European Journal of Combinatorics*, 28(1):202–213, 2007.
- [27] Vít Jelínek and Toufik Mansour. Matchings and partial patterns. *The Electronic Journal of Combinatorics*, 17:R158, 2010.
- [28] Grigorios Koulouras and Martin C. Frith. Significant non-existence of sequences in genomes and proteomes. *Nucleic Acids Research*, 49(6):3139–3155, 2021.
- [29] Catherine L. Lawson, Helen M. Berman, Li Chen, Brinda Vallat, and Craig L. Zirbel. The nucleic acid knowledgebase: a new portal for 3d structural information about nucleic acids. *Nucleic Acids Research*, 52(D1):D245–D254, 2023.
- [30] William A. Lorenz, Yann Ponty, and Peter Clote. Asymptotics of RNA shapes. *Journal of Computational Biology*, 15(1):31–63, 2008.
- [31] Xiang-Jun Lu, Harmen J. Bussemaker, and Wilma K. Olson. DSSR: an integrated software tool for dissecting the spatial structure of RNA. *Nucleic Acids Research*, 43(21):e142, 2015.

- [32] Xiang-Jun Lu and Wilma K. Olson. 3DNA: a software package for the analysis, rebuilding and visualization of three-dimensional nucleic acid structures. *Nucleic Acids Research*, 31(17):5108–5121, 2003.
- [33] Thomas J. Macke, David J. Ecker, Robin R. Gutell, Daniel Gautheret, David A. Case, and Rangarajan Sampath. RNAMotif, an RNA secondary structure definition and search algorithm. *Nucleic Acids Research*, 29(22):4724–4735, 2001.
- [34] Austin Montgomery, Georgios Christos Tsiatsianis, Ioannis Mouratidis, Candace S. Y. Chan, Maria Athanasiou, Anastasios D. Papanastasiou, Verena Kantere, Nikos Syrigos, Ioannis Vathiotis, Konstantinos Syrigos, Nelson S. Yee, and Ilias Georgakopoulos-Soares. Utilizing nullomers in cell-free RNA for early cancer detection. *Cancer Gene Therapy*, 2024.
- [35] OEIS Foundation Inc. The On-Line Encyclopedia of Integer Sequences, 2024. Published electronically at <https://oeis.org>.
- [36] Anton I. Petrov, Craig L. Zirbel, and Neocles B. Leontis. WebFR3D—a server for finding, aligning and analyzing recurrent RNA 3D motifs. *Nucleic Acids Research*, 39(suppl_2):W50–W55, 2011.
- [37] Armando J Pinho, Paulo JSG Ferreira, Sara P Garcia, and João MOS Rodrigues. On finding minimal absent words. *BMC Bioinformatics*, 10(1), 2009.
- [38] Yann Ponty. Ensemble algorithms and analytic combinatorics in RNA bioinformatics and beyond, 2020. HDR, Université Paris-Saclay.
- [39] Michela Quadrini. Structural relation matching: an algorithm to identify structural patterns into RNAs and their interactions. *Journal of Integrative Bioinformatics*, 18(2):111–126, 2021.
- [40] Christian M. Reidys and Rita R. Wang. Shapes of RNA pseudoknot structures. *Journal of Computational Biology*, 17(11):1575–1590, 2010.
- [41] Einar Andreas Rødland. Pseudoknots in RNA secondary structures: representation, enumeration, and prevalence. *Journal of Computational Biology*, 13(6):1197–1213, 2006.
- [42] Michael Sarver, Craig L. Zirbel, Jesse Stombaugh, Ali Mokdad, and Neocles B. Leontis. FR3D: finding local and composite recurrent structural motifs in RNA 3D structures. *Journal of Mathematical Biology*, 56(1–2):215–252, 2007.
- [43] Sandra Smit, Kristian Rother, Jaap Heringa, and Rob Knight. From knotted to nested RNA structures: A variety of computational methods for pseudoknot removal. *RNA*, 14(3):410–416, 2008.
- [44] Vincent Vajnovszki. Generating involutions, derangements, and relatives by ECO. *Discrete Mathematics & Theoretical Computer Science*, 12(1):479, 2010.

- [45] F. H. D. van Batenburg, A. P. Gulyaev, C. W. A. Pleij, J. Ng, and J. Oliehoek. PseudoBase: a database with RNA pseudoknots. *Nucleic Acids Research*, 28(1):201–204, 2000.
- [46] Michael S. Waterman. Secondary structure of single-stranded nucleic acids. *Advances in Mathematics Supplementary Studies*, 1:167–212, 1978.
- [47] Marcin Wojdyr. GEMMI: a library for structural biology. *Journal of Open Source Software*, 7(73):4200, 2022.
- [48] Hua-Ting Yao, Cédric Chauve, Mireille Regnier, and Yann Ponty. Exponentially few RNA structures are designable. In *Proceedings of the 10th ACM International Conference on Bioinformatics, Computational Biology and Health Informatics, BCB’19*, page 289–298. Association for Computing Machinery, 2019.
- [49] Tianshuo Zhou, Wei Yu Tang, David H. Mathews, and Liang Huang. Undesignable RNA structure identification via rival structure generation and structure decomposition. In Jian Ma, editor, *Research in Computational Molecular Biology (RECOMB)*, pages 270–287, Springer, Cham, 2024.