

# METHODS FOR ANALYZING RNA PSEUDOKNOTS VIA CHORD DIAGRAMS AND INTERSECTION GRAPHS

RAYAN IBRAHIM AND ALLISON H. MOORE

**ABSTRACT.** RNA molecules are known to form complex secondary structures including pseudoknots. A systematic framework for the enumeration, classification and prediction of secondary structures is critical to determine the biological significance of the molecular configurations of RNA. Chord diagrams are mathematical objects widely used to represent RNA secondary structures and to analyze structural motifs, however a mathematically rigorous enumeration of pseudoknots remains a challenge. We introduce a method that incorporates a distance-based metric  $\tau$  to analyze the intersection graph of a chord diagram associated with a pseudoknotted structure. In particular, our method formally defines a pseudoknot in terms of a weighted vertex cover of a certain intersection graph constructed from a partition of the chord diagram representing the nucleotide sequence of the RNA molecule. In this graph-theoretic context, we introduce a rigorous algorithm that enumerates pseudoknots, classifies secondary structures, and is sensitive to three-dimensional topological features. We implement our methods in MATLAB and test the algorithm on pseudoknotted structures from the bpRNA-1m database. Our findings confirm that genus is a robust quantifier of pseudoknot complexity.

## 1. INTRODUCTION

Ribonucleic acid (RNA) is a molecule essential to many functions of life, notably gene expression, cellular communication, and the storage and transfer of genetic information. The primary structure of RNA refers to the sequence of its four nitrogenous bases adenine (A), guanine (G), cytosine (C), and uracil (U), attached along a sugar-phosphate backbone [14]. It is well known that RNA molecules fold into a variety of secondary and tertiary structures related to their natural functions via complementary Watson-Crick base pairings and other pairings [36, 14]. Common secondary structure motifs include hairpin loops, stems (i.e. ‘stacks’), bulges, interior loops, multiloops, single-stranded regions, and pseudoknots [22], [12, Figure 1] (see Figure 4). A pseudoknot is a secondary structure motif representing a three-dimensional folding pattern. Pseudoknots were first recognized in the study of the turnip yellow mosaic virus [24, 13], but the term was coined in [31]. The simplest type of pseudoknotted structure is an H-type pseudoknot, formed when nucleobases along the loop of a hairpin bond with nucleobases elsewhere along the sequence [6].

A variety of pseudoknot motifs have been characterized including H-type, K-type, L-type, and M-type motifs [18, 1].

RNA secondary structures may be represented graphically with *chord diagrams*, objects common in enumerative combinatorics and topology. In the representation of an RNA secondary structure, bonded pairs are indicated by arcs (‘chords’) along a line segment or circle. A precursor to a chord diagram representing secondary structures appears in [33] as a connection to predictive models using base-pairing matrices [32]. There, RNA secondary structures were defined as simple planar graphs on a set of  $n$  labeled points such that a path along the  $n$  points represents the primary structure, with other edges representing bonds between bases [33, Definition 2.1]. These planar graphs correspond to crossingless chord diagrams, which have since been studied extensively as models for secondary structures [5, 23, 35, 34, 21, 26]. Pseudoknots occur only when the corresponding chord diagrams contain chords that cross each other. Despite the relative simplicity of a chord diagram, there is no agreed-upon method for quantifying the complexity of RNA pseudoknotting. A naive count of crossings overemphasizes contributions from helical stacking, and different methods for reducing parallel bonds may yield different enumerations of pseudoknots. Moreover, existing methods may ignore some topological features of the 3D conformation. We investigate these discrepancies in Sections 2.2 and 3.1.

Our goal is to construct a mathematically rigorous and topologically robust framework for quantifying pseudoknot complexity in RNA, presented in the familiar language of graph theory and building upon conventions implicitly assumed in the bpRNA method [12] and bpRNA-1m database [11]. The bpRNA-1m database aggregates over 100K RNA secondary structures from seven sources [9, 37, 25, 17, 8, 15, 3]. Internal annotation routines and the enumeration of pseudoknots and other structural motifs are conducted via the algorithmic tool bpRNA [12]. This database provides the primary test case for our graph-theoretic procedures. The strategy presented here verifies the reproducibility of both our methods and those of the bpRNA tool [12], while illustrating discrepancies resulting from topological conformations.

Section 2 develops a mathematical formulation of pseudoknotting using chord diagrams and intersection graphs. The relevant graph theoretical background is reviewed in Section 2.1, and we prove a relationship between vertex cover numbers and the genus of a chord diagram in Theorem 2.6. In Section 2.2, we investigate precedent theories of pseudoknotting and convey these notions into our mathematical framework. Specifically, we use a weighted vertex cover of an intersection graph constructed from a partition of the chord diagram corresponding to an RNA molecule to give a precise enumeration of pseudoknots. In Section 3, we introduce the  *$\tau$ -reduction algorithm*, which systematically reduces the complexity of a chord diagram to quantify pseudoknots in a robust manner that takes into account 3D topological features, including nestings, crossings, and a distance-based threshold  $\tau$ . An explicit

MATLAB implementation of our algorithm is provided in Github [16]. Examples 2.11 and 2.12 in this section demonstrate discrepancies in the existing methodology that are corrected by our current methods. In Section 4 we present the results of our analyses. We report quantities of pseudoknots, improving upon previous methods. In [5] it was shown that these basic motifs constitute the irreducible pseudoknots of genus equal to one (see Section 4.3 for a discussion on genus). Moreover, they posit that the topological genus of a chord diagram provides a classification of RNA secondary structures with pseudoknots. Our analysis in Section 4.3 together with Theorem 2.6 confirms that even with additional topological considerations, genus is a robust classifier of pseudoknot complexity.

## 2. COMBINATORIAL THEORY

**2.1. Chord Diagrams.** A *linear chord diagram*  $D$  is a set of  $n$  points on an oriented line segment together with a (partial) matching of the points. Circular chord diagrams are obtained by joining the endpoints of the segment; however we will restrict to linear chord diagrams throughout. We denote chords as pairs  $c = (\ell, r)$ , where  $\ell$  and  $r$  denote left and right endpoints, respectively. Because chord diagrams are matchings,  $\ell < r$ , and no two chords share an endpoint. By convention, a set of chords is indexed by left endpoints.

**Definition 2.1.** For any two chords  $c_1$  and  $c_2$ , there are three possibilities:

- (a)  $c_1$  and  $c_2$  form a *crossing*:  $\ell_1 < \ell_2 < r_1 < r_2$ .
- (b)  $c_1$  and  $c_2$  form a *nesting*:  $\ell_1 < \ell_2 < r_2 < r_1$ .
- (c)  $c_1$  and  $c_2$  are *independent*:  $\ell_1 < r_1 < \ell_2 < r_2$ .

Further, a  $k$ -crossing is a set of chords  $(\ell_1, r_1), (\ell_2, r_2), \dots, (\ell_k, r_k)$  such that  $\ell_1 < \ell_2 < \dots < \ell_k < r_1 < r_2 < \dots < r_k$ . A  $k$ -nesting is a set of chords  $(\ell_1, r_1), (\ell_2, r_2), \dots, (\ell_k, r_k)$  such that  $\ell_1 < \ell_2 < \dots < \ell_k < r_k < r_{k-1} < \dots < r_1$ .

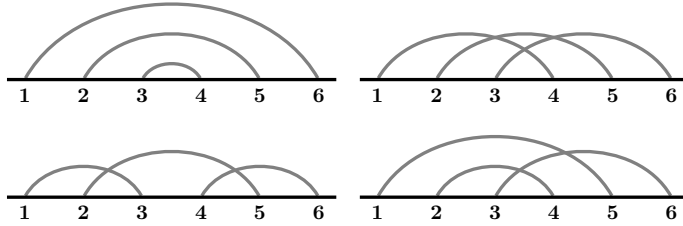


FIGURE 1. Top left: A 3-nesting. Top right: A 3-crossing. Bottom left: Two 2-crossings. Bottom right: A 2-nesting and two 2-crossings.

The following definitions will become useful in formalizing pseudoknots in chord diagrams. We use the notation  $(a, b)$  for the open interval between  $a$  and  $b$ .

**Definition 2.2** (Chord Obstructed). Let  $c$  and  $c'$  be chords and let  $U = (\ell, \ell') \cup (r, r')$ , i.e.  $U$  can be thought of as the set of bases between the left endpoints and right endpoints of  $c$  and  $c'$ . We say  $c$  and  $c'$  are *chord obstructed* if there is some chord  $c''$  such that either  $\ell'' \in U$  or  $r'' \in U$ . A set of three or more chords  $S = \{c_1, c_2, \dots, c_k\}$  is chord obstructed if consecutive chords  $c_i$  and  $c_{i+1}$  are chord obstructed for some  $i = 1, 2, \dots, k-1$ .

**Definition 2.3** (Segment). A *segment* of a chord diagram  $D$  is a maximal nonempty set of chords  $S = \{c_1, c_2, \dots, c_k\}$  forming a  $k$ -nesting such that  $S$  is not chord obstructed.

Note that the set of segments  $\mathcal{S}$  partitions the set of chords  $C$ . In crossingless chord diagrams, there is a natural poset structure on the set of segments defined by  $S \prec S'$  if  $S'$  is nested in  $S$ .

The *intersection graph*  $G$  of a chord diagram  $D$  is the graph whose vertices are the chords of  $D$ , and such that two vertices in  $G$  are adjacent if their corresponding chords in  $D$  form a crossing. The use of intersection graphs is well established, see for example [19]. Variations on the concept of an intersection graph arise numerous times in the biology literature under different names. For example, in [18], there is the notion of a conflict graph, whose vertex set comprises helices in the RNA structure and edges signify the crossing of chords corresponding to the helices. In [29] the concept of an element-contact graph is introduced, in particular the stem-loop-contact graph (SLCG) [29, Figure 7]. The segment graph of the bpRNA database [12] is another such example.

We will investigate the following graph theoretic invariants with respect to RNA structures in Section 4. An *independent set* is a set of vertices such that no two vertices in the set are adjacent. The maximum *clique number*  $\omega(G)$  is the maximum size of a complete subgraph of  $G$ . A *vertex cover* of a graph  $G$  is a set  $A \subseteq V(G)$  such that for every edge  $xy \in E(G)$  either  $x \in A$  or  $y \in A$ . The *vertex cover number* of a graph  $G$ , denoted  $\beta(G)$ , is the number of vertices in a minimum vertex cover. The *weight* of a vertex cover  $A$  in a vertex-weighted graph  $G$  is  $\sum_{v \in A} w(v)$ . Note that a *minimum weight vertex cover* of a weighted graph  $G$  is not necessarily a minimum cardinality vertex cover. (Consider for example the path graph  $P_3$  with weights 1,5,3).

If  $C$  is a vertex cover of  $G$ , then the graph  $G - C$  contains no edges, as by definition every edge of  $G$  must have an endpoint in  $C$ . Similarly, if  $I$  is an independent set of  $G$ , then every edge of  $G$  has at least one end point in  $G - I$ . Thus we have the following observation.

**Observation 2.4.** *Let  $G$  be a graph and let  $I$  and  $C$  be an independent set and vertex cover respectively. Then  $V(G) \setminus I$  is a vertex cover, and  $V(G) \setminus C$  is an independent*

set. Moreover, the complement of a minimum weight vertex cover is a maximum weight independent set.

**Definition 2.5** (Genus). The *genus* of a chord diagram  $D$ , denoted  $\gamma(D)$ , is half of the rank of the adjacency matrix of the intersection graph with  $\mathbf{Z}_2$  coefficients [20].

Equivalently, the genus of a chord diagram is equal to the topological genus of the surface obtained by regarding the chord diagram as a band surgery diagram [5]. One way to calculate the genus  $g$  of a chord diagram  $D$  is via the formula

$$g = \frac{P - L}{2}$$

where  $P$  is the number of chords (or base pairs) and  $L$  is the number of closed loops in the corresponding double-line diagram [5] (see Figure 2.) For a graph  $G$  we define  $\text{rank}(G)$  (resp.  $\text{rank}_p(G)$ ) to be the rank of the adjacency matrix of  $G$  with coefficients in  $\mathbb{R}$  (resp. coefficients in  $\mathbb{F}_p$ ).

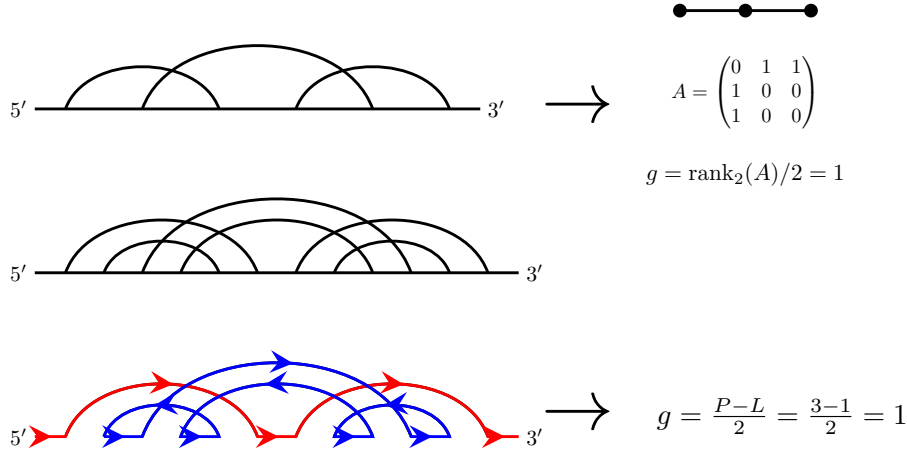


FIGURE 2. Two methods of calculating the genus of a given chord diagram.

**Theorem 2.6.** Let  $D$  be a chord diagram and let  $G$  be the intersection graph of  $D$ . If  $G$  is acyclic, then  $\gamma(D) = \beta(G)$ .

*Proof.* Because  $\beta(G)$  and  $\gamma(G)$  are additive over disjoint unions, without loss of generality we may assume  $G$  is a tree of order  $n$ . The statement is easily verified for  $n \leq 4$ . We will use induction on  $n$ . Let  $T$  be a tree on  $n \geq 5$  vertices with adjacency matrix  $A$ . If  $T$  has a vertex  $v$  with at least two leaf neighbors, then we remove a leaf  $\ell$  adjacent to  $v$  to form  $T'$ . By the induction hypothesis,  $\text{rank}_2(T')/2 = \beta(T')$ . Adding back  $\ell$  to form  $T$ , we have  $\text{rank}_2(T') = \text{rank}_2(T)$  as the row corresponding to  $\ell$  in  $A$  is identical to the rows corresponding to the other leaf neighbors of  $v$ . Note

in a graph, for any vertex with a leaf neighbor, there is a minimum vertex cover containing that vertex. Thus, we have  $\beta(T) = \beta(T')$ .

Now we assume every vertex of  $T$  has at most one leaf as a neighbor. Then there is some leaf  $\ell$  with a neighbor  $v$  of degree two; to identify such a vertex  $v$ , find a maximum path and let  $v$  be adjacent to a leaf. Let  $\{\ell, w\} = N(v)$ . Let  $T' = T - \{v, \ell\}$ . By the induction hypothesis,  $\gamma(T') = \beta(T')$ .

First we claim  $\beta(T) = \beta(T') + 1$ . Indeed, no minimum vertex cover of  $T'$  will cover the edge  $v\ell$ , so  $T$  requires one more vertex in addition to a minimum vertex cover of  $T'$  to cover all edges of  $T$ .

Next we claim that  $\text{rank}_2(T) = \text{rank}_2(T') + 2$ . The adjacency matrix of  $T$  is given by

$$A = \begin{array}{c} \begin{array}{cccc|ccc} & & & & w & v & \ell \\ & & & & * & 0 & 0 \\ & & & & \vdots & \vdots & \vdots \\ & & & & * & 0 & 0 \\ w & * & \cdots & * & 0 & 1 & 0 \\ v & 0 & \cdots & 0 & 1 & 0 & 1 \\ \ell & 0 & \cdots & 0 & 0 & 1 & 0 \end{array} \end{array}$$

where the upper left block corresponds to the adjacency matrix  $A'$  of  $T'$ . Applying the row operation  $\text{row } w - \text{row } \ell = \text{row } v$  and column operation  $\text{col } w - \text{col } \ell = \text{col } w$  we obtain the matrix

$$B = \begin{array}{c} \begin{array}{cccc|ccc} & & & & w & v & \ell \\ & & & & * & 0 & 0 \\ & & & & \vdots & \vdots & \vdots \\ & & & & * & 0 & 0 \\ w & * & \cdots & * & 0 & 0 & 0 \\ v & 0 & \cdots & 0 & 0 & 0 & 1 \\ \ell & 0 & \cdots & 0 & 0 & 1 & 0 \end{array} \end{array}.$$

We have that  $\text{rank}_2(A) = \text{rank}_2(B) = \text{rank}_2(A') + 2$ , therefore  $\gamma(T) = \gamma(T') + 1$ . This completes the induction.  $\square$

**Remark 2.7.** In general, for acyclic graphs  $\gamma(G) \neq \beta(G)$ . In particular, for complete graphs, one may observe that  $\beta(K_n) - \gamma(K_n) = k - \lfloor \frac{k}{2} \rfloor - 1$ .

Finally, we remark that the genus of a chord diagram  $D$  containing an  $r$ -nesting  $C = \{c_1, \dots, c_r\}$  that is not chord obstructed is equal to the genus of the diagram with  $r - 1$  chords of  $C$  removed, i.e.  $D - \{c_2, \dots, c_r\}$ . The rank of a matrix can be thought of as the maximum number of linearly independent rows or columns in the matrix. It can be seen from the adjacency matrix  $A$  of the intersection graph that  $c_2, \dots, c_r$  correspond to identical rows in  $A$ , and are thus linearly dependent.

**2.2. Pseudoknotted Structures.** Analyzing the role of pseudoknots in various RNA processes motivates the study of pseudoknot complexity for comparison, characterization of motifs, and prediction of RNA secondary structure [12, 5, 23, 30]. Practical definitions of a pseudoknot vary in the biological literature [5, 23]. Because we will ultimately be interested in quantifying pseudoknots aggregated by the bpRNA-1m database [12], we start from definitions presented there. In this database, a pseudoknotted structure is characterized as having base pair positions that cross in the sense of Definition 2.1(a); the working definition of a ‘pseudoknot base pair’ is one belonging to a minimal set of base pairs that results in a pseudoknot-free structure once removed. There is some ambiguity in these concepts; for example, a pair of kissing hairpins with intersection graph weighted  $[a, a + b, b]$  demonstrates that such a minimal set is not unique. Further ambiguities may result from helices formed in secondary structures (see Section 2.4). By default, the number of pseudoknots in any given secondary structure reported in [12] is the number returned by an algorithm that identifies some minimal set of pseudoknot base pairs. We review their algorithm, and translate corresponding notions of pseudoknotting into the language of chord diagrams, as follows.

To make the notion of a pseudoknot, and more specifically the annotation of multiple pseudoknots, more precise, [12] introduces the notion of a segment of RNA secondary structure. An RNA segment is described as a region of duplexed RNA, possibly containing bulges or internal loops. In combinatorial terms, RNA segments correspond to the segments  $s \in \mathcal{S}$  of the chord diagram  $D$  representing the secondary structure, as in Definition 2.3. The segments partition the set of chords  $C$ , and ordering the base sequence from the 5'-end to 3'-end indexes each segment by its leftmost endpoint. In [12], the following is observed; we provide a restatement in terms of chord diagrams.

**Theorem 2.8** ([12]). *Let  $\mathcal{S}$  be the segment partition of a linear chord diagram  $D$  and let  $S, S' \in \mathcal{S}$ . If there are chords  $c \in S$  and  $c' \in S'$  such that  $c$  and  $c'$  are crossed, then any pair of chords from  $S$  and  $S'$  cross.*

*Proof.* Recall that a segment of size  $k$  is a maximal  $k$ -nesting in which pairs of consecutive chords are not chord obstructed. Let  $\mathcal{S}$  be a segment partition of  $D$  and let  $S, S' \in \mathcal{S}$  where  $S < S'$  in the indexing of segments by left endpoints. Let  $c \in S$  and  $c' \in S'$  such that  $c$  and  $c'$  cross. Let  $\ell_{\max}$  be the maximum left endpoint of  $S$  and  $r_{\min}$  and  $r_{\max}$  be the minimum and maximum right endpoints of  $S$  respectively.

By the left endpoint indexing, all left endpoints of chords in  $S'$  must be greater than  $\ell_{\max}$ . If some left endpoints of  $S'$  are less than  $r_{\min}$  and some greater, then  $S'$  is chord obstructed by the right endpoints of  $S$ . If all left endpoints of  $S$  are greater than  $r_{\min}$ , either  $S$  is chord obstructed or no chords of  $S$  and  $S'$  cross. Thus, all left endpoints of  $S'$  must lie between  $\ell_{\min}$  and  $r_{\min}$ . Since  $c$  and  $c'$  cross, and  $S$  is not

chord obstructed, it must be that  $r' > r_{\max}$ , i.e.  $c'$  crosses every chord in  $S$ . No other right endpoint of  $S'$  is less than  $r_{\max}$ , as otherwise  $S$  or  $S'$  are chord obstructed. Thus any pair of chords from  $S$  and  $S'$  cross.  $\square$

In other words, Theorem 2.8 says that two segments  $s, s' \in \mathcal{S}$  cross whenever any chords  $c \in s, c' \in s'$  cross. The *segment graph*  $G_{\mathcal{S}}$  of a chord diagram is a variation on an intersection graph whose vertex set is the set of segments, where two vertices are adjacent if their segments cross [12]. Vertices in  $G_{\mathcal{S}}$  are weighted by the number of chords contained in their corresponding segments. The terminology *PK-segment* refers to any segment  $s \in \mathcal{S}$  that crosses any other segment. Let  $H_{\mathcal{S}} \subset G_{\mathcal{S}}$  be the subgraph with isolated vertices removed, referred to as the *PK-segment graph* in [12].

In [12], pseudoknotted structures are identified by finding a maximum weight independent set  $I$  in  $H_{\mathcal{S}}$  via a heuristic approach, with an exact algorithm used in the specific case of components which are paths. As we have formalized above, the set  $P = V(G_{\mathcal{S}}) - I$  is a minimum weight vertex cover of  $G_{\mathcal{S}}$ . That is,  $P$  is a set of segments of minimum cardinality in  $\mathcal{C}$  that when removed from  $D$  leave a pseudoknot-free structure.

We may now formally quantify the size of a pseudoknotted structure according to the conventions of the bpRNA-1m database, as implied by the algorithms of [12].

**Definition 2.9.** (Pseudoknotted Structures - *bpRNA Segment Graph Method*) A secondary structure is *pseudoknotted* if its segment graph contains at least one edge and is called *pseudoknot-free* otherwise. The *number of pseudoknots* in a pseudoknotted structure is the minimum cardinality over all vertex covers of minimum weight of the segment graph.

Any segment contained in a minimum-cardinality minimum-weight vertex cover may be called simply ‘a pseudoknot.’ It is important to note that in [12], the number of pseudoknots is not the number of chords in the corresponding cover, in general.

**Example 2.10.** Figure 3 shows the chord diagram representing the secondary structure of tRNA (76-MER) found in *Escherichia coli* with 8 segments. The PK-segments are  $\{2, 3, 4, 5, 6, 8\}$ . The maximum weight independent set in  $H_{\mathcal{S}}$  is  $\{3, 8\}$ . The minimum weight vertex cover in  $G_{\mathcal{S}}$  is  $\{2, 4, 5, 6\}$ , indicating four pseudoknots in this structure according to the conventions of [12].

**2.3. Crossingless Secondary Structures.** Let  $D$  be a chord diagram with no crossings, and let  $\mathcal{S}$  be the segment partition of  $D$ . We say an unpaired base is nested in a segment  $S$  if it is between the left and right endpoints of the innermost chord of  $S$ . A *stem* in  $D$  is a  $k$ -nesting with no other bases between the endpoints of any two consecutive chords. Segments are composed of stems. Let  $S$  be a segment and let  $c$  and  $c'$  be two consecutive chords in a segment  $S$ , with  $\ell < \ell'$ . If there



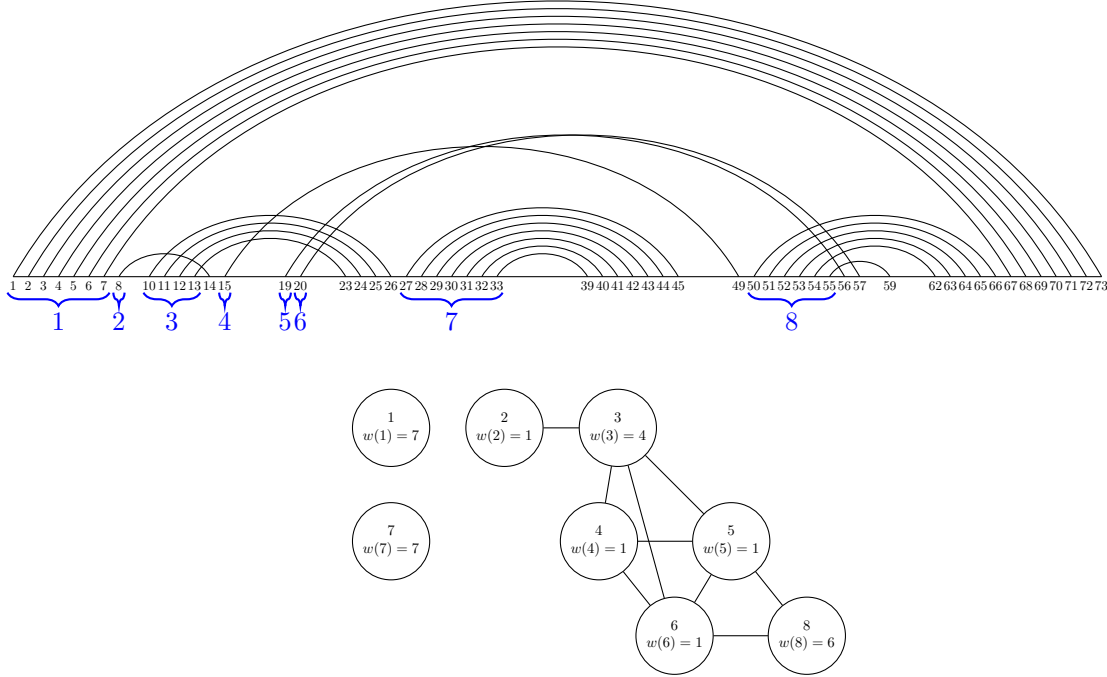


FIGURE 3. (Top) Linear chord diagram of transfer RNA molecule of type 76-MER from *Escherichia coli* (PDB.652) with segments labeled. (Bottom) Segment graph  $G_S$  with weights.

are sequences of unpaired bases in the intervals  $(\ell, \ell')$  and  $(r, r')$ , i.e.  $\ell' - \ell \geq 2$  and  $r - r' \geq 2$ , then those two sequences together comprise an *interior loop*. If exactly one of  $(\ell, \ell')$  or  $(r, r')$  contains a sequence of unpaired bases, that sequence is a *bulge*. If two or more segments  $\mathcal{T}$  are nested in  $S$ , then there is a *multiloop* composed of all unpaired bases  $b$  nested in  $S$  and not nested in any segment in  $\mathcal{T}$ . Note multiloops may have length zero (that is, a multiloop may be the empty set). The *exterior loop* of  $D$  is the set of all unpaired bases which are not nested in any chord. If  $b_0$  is the first paired base and  $b_f$  is the last paired base in the base sequence, the set of unpaired bases less than  $b_0$  and greater than  $b_f$  are a part of the exterior loop called the *dangling ends*. The exterior loop can be thought of as the multiloop gained from an imaginary base pair bonding of the 5'-end and 3'-end, under which all chords are nested. The crossingless secondary structures are illustrated in Figure 4.

**2.4. Qualitatively Similar Pseudoknotted Structures.** Identifying nested chords in a chord diagram is a common strategy for reducing the complexity of the combinatorial analysis of secondary structures because it allows for the reduction of helices to a single chord, or alternatively to a single vertex in an intersection graph. For

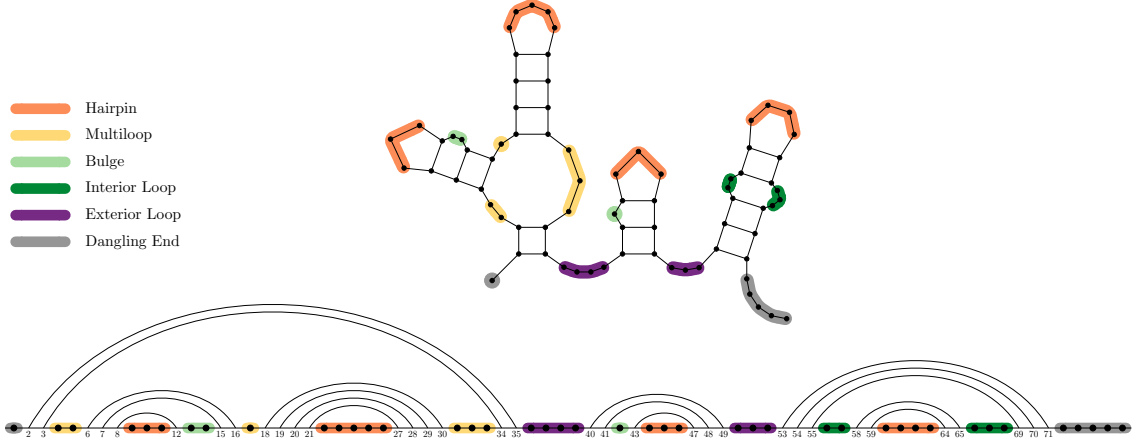


FIGURE 4. An illustrated example of secondary structures in a crossless chord diagram. There are eight stems in this structure

example, nested chords are identified as segments in [12] (Definition 2.3 above), as ‘stacks’ in [5], as ‘shadows’ in [23], or as edges in tree diagrams of pseudoknot-free structures in [2] and elsewhere. Such reductions also preserve some invariants of interest, for example the topological genus (see Definition 2.5 and [5]). In contrast to  $r$ -nestings,  $r$ -crossings in chord diagrams are typically left untouched by simplification algorithms. Consequently, the existence of an  $r$ -crossing implies the existence of at least  $r$  pseudoknots according to the conventions of [12].

Consider the following two examples.

**Example 2.11** (Discrepancies due to weights). Consider two weighted segment graphs, each isomorphic to  $P_3$ , with weights  $(1, 5, 3)$  and  $(1, 5, 4)$ , respectively. Such graphs represent nearly identical  $K$ -type secondary structures which differ only by a single bonded pair in the third stem. By [12] and Definition 2.9, the minimum cardinality over vertex covers of minimum weight determines the number of pseudoknots: 2 and 1, respectively. Note that the addition of one pair results in a *decrease* in the number of pseudoknots.

**Example 2.12** (Discrepancies due to  $r$ -crossings). Consider a nucleotide sequence that contains a repetitious subsequence appearing in reverse. For an example, we follow Figure 5. This structure contains subsequence  $\sigma = X'Y'Z'$ , complementary sequence  $\sigma' = ZYX$ , and reverse complementary sequence  $\bar{\sigma}' = XYZ$ . Bonds formed between  $\sigma$  and  $\sigma'$  result in an  $r$ -nesting (Figure 5(A)), whereas bonds formed between  $\sigma$  and  $\bar{\sigma}'$  form an  $r$ -crossing (Figure 5(B)). We assume here that either set of bonds is possible. In particular, the parity of number of half-turns in the helical stem may determine whether  $\sigma'$  or  $\bar{\sigma}'$  is nearer to  $\sigma$  in a 3D conformation. Despite the

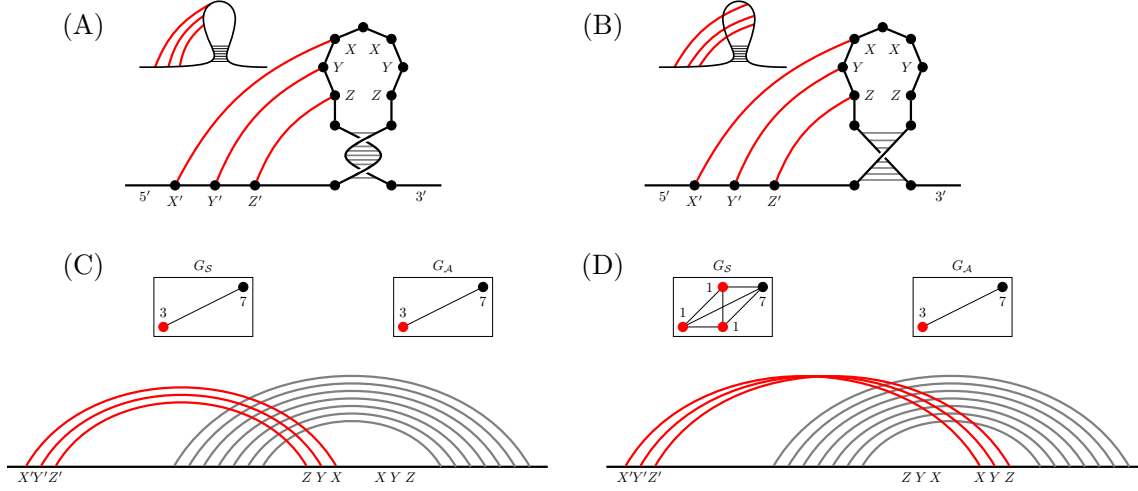


FIGURE 5. Two closely related instances of bonding between two hairpins

heuristic similarity of the resulting pseudoknotted structures, their segment graphs differ significantly. By the conventions of [12] (Definition 2.9) the two structures are of pseudoknotting size 1 and 3, respectively. Moreover, the difference increases with additional base pairing in the hairpin stem.

The general observation is that variances in bonding from spatial conformations (e.g. helical twisting) may result in a quantification of pseudoknotting that is artificially high. Of 30 structures exhibiting the most pseudoknotting in Table 1, there are 28 that contain both a complementary and reverse complementary sequence (possibly not contiguous) near the site of the 3-crossing. One example RNA structure involving a 3-crossing and 5-nesting is bpRNA\_CRW\_55315. This structure contains a 3-crossing with left bases GCA at indices 2107, 2108, and 2112 and right bases AGU at indices 2164, 2165, and 2167, which is the reverse of the triple UGA at indices 2162, 2163, and 2164.

### 3. METHODS

To further reduce complexity and to more effectively relate similar secondary structures, we propose in this section an alternative method for quantifying the size of a pseudoknotted structure and a new simplification algorithm that identifies both  $r$ -crossings and  $r$ -nestings in chord diagrams. In addition to handling complexity issues arising from  $r$ -crossings, this method will also eliminate some discrepancies resulting from weighted graphs and include a parameter accounting for distance between nucleotides.

**3.1. The  $\tau$ -Segment Graph Method.** In this subsection, we give a partitioning procedure that identifies  $r$ -crossings in a manner similar to that of  $r$ -nestings. We implement the procedure in an algorithm that incorporates an additional distance parameter  $\tau$  in terms of the nucleotide sequence.

An *augmented segment* of a chord diagram is a maximal nonempty set of chords  $S = \{c_1, c_2, \dots, c_k\}$  forming a  $k$ -nesting or a  $k$ -crossing which is not chord obstructed. Revisiting Examples 2.11 and 2.12 above, notice that an intersection graph  $G_{\mathcal{A}}$  produced with augmented segments would yield two pseudoknotted structures of size 1 in Example 2.11 and two structures of size 2 in Example 2.12.

The *chord distance* between two chords  $c_1 = (\ell_1, r_1)$  and  $c_2 = (\ell_2, r_2)$  is

$$d(c_1, c_2) = \max\{|\ell_1 - \ell_2|, |r_1 - r_2|\}.$$

We say that a pair of nested or crossed chords  $c_1, c_2$  are  $\tau$ -near if  $d(c_1, c_2) \leq \tau$  and  $c_1$  and  $c_2$  are not chord obstructed. A  $k$ -crossing or  $k$ -nesting  $\{c_1, \dots, c_k\}$  is  $\tau$ -near if for each  $i = 1, \dots, k-1$  we have  $c_i$  and  $c_{i+1}$  are  $\tau$ -near. A  $\tau$ -segment of a chord diagram is a maximal nonempty set of chords  $S = \{c_1, c_2, \dots, c_k\}$  forming a  $\tau$ -near  $k$ -nesting or a  $\tau$ -near  $k$ -crossing. As with segments or augmented segments, the set of  $\tau$ -segments  $\mathcal{S}_\tau$  also partitions the set of chords  $C$ . We have the following statement.

**Theorem 3.1.** *Let  $\mathcal{S}_\tau$  be the  $\tau$ -segment partition of a linear chord diagram  $D$  and let  $S_1, S_2 \in \mathcal{S}_\tau$ . If there are chords  $c \in S_1$  and  $c' \in S_2$  such that  $c$  and  $c'$  are crossed, then for every pair  $c \in S_1$  and  $c' \in S_2$ , the chords  $c$  and  $c'$  are crossed.*

*Proof.* The proof is analogous to that of Theorem 2.8. □

Generalizing the segment graph, we may now define the  $\tau$ -segment intersection graph  $G_\tau$  to be the weighted graph whose vertex set is the set of  $\tau$ -segments, where two vertices are adjacent if the  $\tau$ -segments cross. When  $\tau = 0$ , define  $G_{\tau=0} := G_{\mathcal{S}}$ . When  $\tau = \infty$ , the graph  $G_{\mathcal{A}} := G_\infty$  is the intersection graph of the augmented segment partition. The notation  $D_\tau$  will indicate the chord diagram corresponding to  $G_\tau$ , in which each  $\tau$ -segment corresponds to a chord. As with  $G_{\mathcal{S}}$  and  $G_{\mathcal{A}}$  we use the notation  $D_{\mathcal{S}}$  and  $D_{\mathcal{A}}$  analogously. We may now formally revise the method for quantifying the size of pseudoknotted structures in RNA.

**Definition 3.2.** (Pseudoknotted Structures -  $\tau$ -Segment Graph Method) A secondary structure is  $\tau$ -pseudoknotted if  $G_\tau$  contains at least one edge and is called *pseudoknot-free* otherwise. For  $\tau \geq 1$ , the *number of pseudoknots* is the minimum cardinality of a vertex cover of  $G_\tau$ . For  $\tau = 0$ , the number of pseudoknots is the minimum cardinality over all vertex covers of minimum weight of the segment graph.

The definition in the case of  $\tau = 0$  is explicitly made to agree with the conventions of [12].

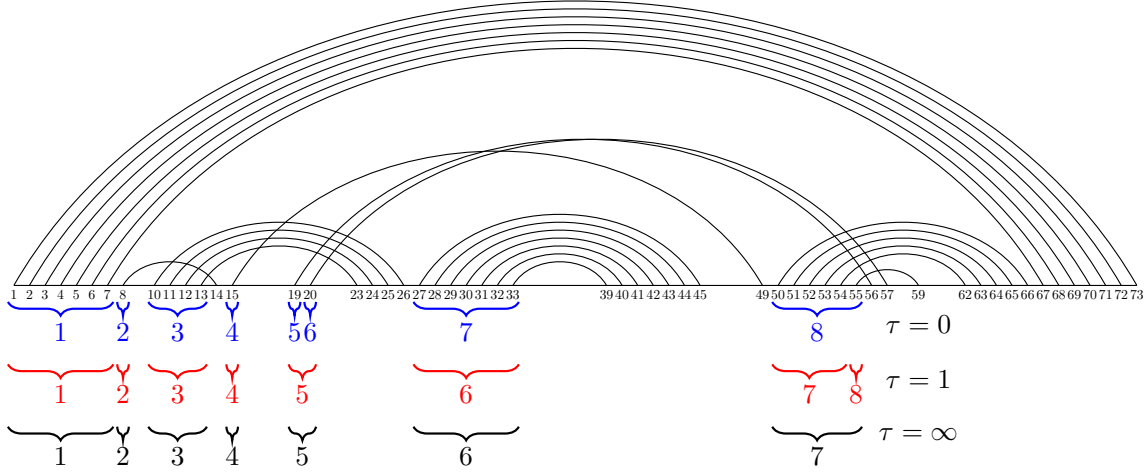


FIGURE 6. Different tau-partitions of bpRNA\_PDB\_652 for  $\tau = 0, 1$ , and  $\infty$ .

**Example 3.3.** Figure 6. The maximum distance between any two non-chord obstructed chords is three. That is, the  $\tau$ -segment partition is the same for all  $\tau \geq 3$ . This RNA structure has four pseudoknots according to the conventions of [12], and two pseudoknots according to Definition 3.2. The main difference comes from chords (19, 56) and (20, 57) becoming part of the same segment for  $\tau$  large enough.

Algorithm 1 in the next section implements the  $\tau$ -segment partition procedure. The input to the algorithm is a chord diagram with set of chords  $C$  (a list of base pairs indexed by left endpoint) and a non-negative integer parameter  $\tau$ . Selecting a pair of chords  $c, c' \in C$ , the algorithm loops to build the  $\tau$ -segment containing  $c$ . If  $\tau = 0$ , it checks whether the  $c, c'$  are nested and not chord obstructed. If  $\tau > 0$ , the algorithm checks whether  $c$  and  $c'$  are  $\tau$ -near and not chord obstructed. If the criterion is met,  $c'$  is added to the segment containing  $c$  and the next pair is selected. If not, the segment is closed. The next unvisited pair of chords are then selected and the process repeats to build the next segment until all chords have been exhausted.

The Github repository [16] contains Algorithm 1 implemented in MATLAB in the function called `findSegments`.

**3.2. Pseudoknot Quantification Process.** Here, we apply the definitions and algorithm above to the bpRNA-1m(90) database [12]. The database bpRNA-1m(90) is a subset of bpRNA-1m restricted to the 28,370 RNA secondary structures with less than 90% sequence similarity. This database contains 3,320 RNA structures reported to contain at least one pseudoknot, i.e. structures whose segment graphs contain at least one edge, with a total of 7,164 pseudoknots reported according to

**Algorithm 1** Tau-Segment Partition

---

**Require:**  $C = \{c_1, c_2, \dots, c_k\} = \{(\ell_1, r_1), (\ell_2, r_2), \dots, (\ell_k, r_k)\}$ ,  $\tau$

**Ensure:**  $\mathcal{C}_\tau$

```

 $S \leftarrow \{(\ell_1, r_1)\}$ 
if  $\tau > 0$  then
  for  $1 \leq i \leq k - 1$  do
     $d \leftarrow \max\{|\ell_i - \ell_{i+1}|, |r_i - r_{i+1}|\}$  ▷ Distance.
    if  $d \leq \tau$  and  $\neg \text{isChordObstructed}(c_i, c_{i+1})$  then
       $S \leftarrow \text{append}(S, (\ell_{i+1}, r_{i+1}))$ 
    else
       $\mathcal{C}_\tau \leftarrow \text{append}(\mathcal{C}_\tau, S)$  ▷ Store Segment
       $S \leftarrow \{c_{i+1}\}$  ▷ Initialize new segment
    end if
  end for
else
  for  $1 \leq i \leq k - 1$  do
    if  $\text{isNested}(c_i, c_{i+1})$  and  $\neg \text{ChordObstructed}(c_i, c_{i+1})$  then
       $S \leftarrow \text{append}(S, c_{i+1})$ 
    else
       $\mathcal{C}_\tau \leftarrow \text{append}(\mathcal{C}_\tau, S)$  ▷ Store Segment
       $S \leftarrow \{c_{i+1}\}$  ▷ Initialize new segment
    end if
  end for
end if
 $\mathcal{C}_\tau \leftarrow \text{append}(\mathcal{C}_\tau, S)$  ▷ Account for final segment.

```

---

Definition 2.9 (the prior bpRNA Segment Graph Method). To analyze the data, we implement both the segment and  $\tau$ -segment graph methods and analyze secondary structures via MATLAB code [16].

Chord diagrams associated with RNA structures are stored as MATLAB arrays. The input to the Tau-Segment Partition algorithm is the set of chords  $C$  from a chord diagram  $D$  and an integer parameter  $\tau$ . To carry out any of the above methods, we first call the **findSegments** function to create a segment partition of the chord diagram. The parameter  $\tau$  determines which segment partition is created, with  $\tau = 0, \infty, 1 < \tau < \infty$  corresponding to the segment partition, augmented segment partition, and  $\tau$ -partition, respectively. See Figure 6.

Depending on whether the segment partition contains any segments which cross each other, one of two subroutines is implemented. If the segment partition contains no segments which cross, then the secondary structures of the chord diagram are analyzed by the function **classifyBases**. This function outputs the primary base

sequence with each base classified as belonging to one of the secondary structure motifs given in Section 2.3.

If the segment partition contains segments which cross, then we construct the corresponding weighted segment graph with the function `makeSegmentGraph`. The pseudoknots of the chord diagram are then identified and analyzed by the `findPKs` process:

- (1) Find list of all maximal independent sets  $I$ .
- (2) If  $\tau = 0$ :
  - (a) Calculate weight of each set  $I \in \mathcal{I}$ .
  - (b) Subset all maximum weight sets  $\mathcal{I}' \subseteq \mathcal{I}$ .
  - (c) Subset maximum cardinality maximum weight sets  $\mathcal{I}'' \subseteq \mathcal{I}' \subseteq \mathcal{I}$ .
- (3) If  $\tau > 0$ :
  - (a) Subset all maximum cardinality sets  $\mathcal{I}' \subseteq \mathcal{I}$ .
  - (b) Calculate weight of each set  $I \in \mathcal{I}'$ .
  - (c) Subset maximum weight maximum cardinality sets  $\mathcal{I}'' \subseteq \mathcal{I}' \subseteq \mathcal{I}$ .
- (4) Dualize the independent sets  $\mathcal{I}''$  to vertex covers  $\mathcal{P}$ .
- (5) Select first vertex cover  $P \in \mathcal{P}$  with respect to the lexicographical ordering from the indexing of segments by left endpoints.

Step (1) applies the Bron-Kerbosch algorithm [7, 4] to find all maximal cardinality independent sets. The importance of the lexicographical ordering in step (5) will become apparent after Example 3.4 below.

The output of this process is a vertex cover  $P$ . In the case that  $\tau = 0$ , this vertex cover represents a pseudoknotted structure by Definition 2.9, where the number of pseudoknots is quantified by the minimum cardinalities of vertex covers of minimum weight. In the cases where  $\tau > 0$ , the minimum-weight minimum-cardinality vertex cover represents a  $\tau$ -pseudoknotted structure by Definition 3.2.

After identifying pseudoknots, one may still want to analyze the crossingless secondary structures of the RNA molecule. Therefore the last part of the process removes all chords which compose a pseudoknot, leaving a set of chords  $C' = C - P$ . In the case of  $\tau = 0$ , the set of chords  $C'$  is crossingless. In the case  $\tau > 0$ , it may be that chords in  $C'$  cross, however  $C'$  comprises independent crossings and nestings, i.e. no two segments cross in the augmented segment partition of  $C'$ . Setting  $\tau = \infty$  (which corresponds to an augmented segment partition), we enter the secondary structure classification subroutine (`classifyBases`) assessing secondary structure using the definitions given in Section 2.3, with independent crossings handled as nestings (and thus a type of stem).

This entire process is summarized in a flow chart (see Figure 8) in the Appendix.

**Example 3.4** (Lexicographical Ordering Matters). Here we apply the  $\tau = 0$  reduction method to the RNA structure 3DIG from the Protein Data Bank [28, 27] (see

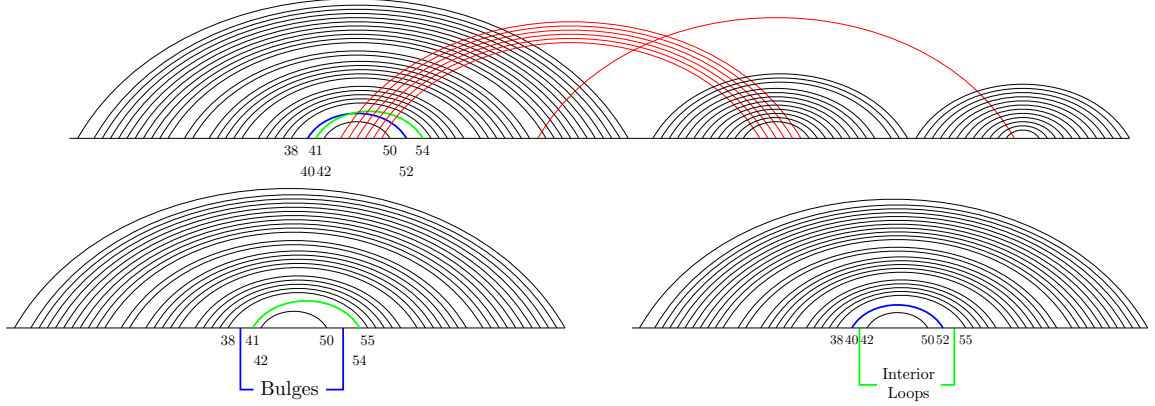


FIGURE 7. RNA structure with bpRNA reference name PDB.455 and PDB reference name 3DIG. Structures such as nestings not relevant for the discussion have been omitted for clarity.

Figure 7). There are two choices for a minimum-cardinality minimum-weight vertex cover, here of cardinality three. Namely, both covers contain the two segments highlighted in red, and differ by whether the cover contains segment  $\{(40, 52)\}$  or segment  $\{(41, 54)\}$ . Removing either cover yields a chord diagram with no crossings, and finding the secondary structures, we may obtain two different pseudoknot types depending on the cover removed. If we remove the cover containing segment  $\{(40, 52)\}$ , then the pseudoknot corresponding to  $\{(40, 52)\}$  connects a bulge to a bulge. However if instead we remove the cover containing segment  $\{(41, 54)\}$  then the pseudoknot corresponding to  $\{(41, 54)\}$  connects an interior loop to another interior loop. As a result, the two choices for a vertex cover have a different effect on the secondary structure classification and consequently pseudoknot typing.

#### 4. DISCUSSION

We applied the  $\tau$ -Segment Graph Method with  $\tau = 0$  to independently verify the quantities reported in bpRNA-1m(90) [12]. By Definition 2.9, structures are pseudoknotted when their segment graphs contain at least one edge, and the number of pseudoknots is quantified by the minimum cardinalities of vertex covers of minimum weight. In agreement with [12], we obtained 3,320 graphs containing at least one edge in  $G_S$  from RNA structures and a total quantity of 7,164 pseudoknots, as determined by the sum over the cardinalities of the vertex covers. Applying the  $\tau$ -Segment Graph Method with  $\tau = \infty$  (the augmented segment graph method) for every structure in bpRNA-1m(90), we found a minimum vertex cover for each structure. The number of pseudoknots (the sum of vertex cover numbers over all graphs) was 6,548 with this method.



With  $\tau = 0$ , we found 31 unique RNA structures containing 13 or more pseudoknots. These structures are listed in Table 1. When the method was applied with  $\tau = \infty$ , we found that the same 31 structures contained the most pseudoknots amongst all structures in the database. With the exception of the last structure, *Oceanobacillus iheyensis*, from  $\tau = 0$  to  $\tau = \infty$  there was a uniform decrease in pseudoknotting by two that resulted from a single 3-crossing being consolidated into one segment by  $\tau$ -reduction. This uniformity in behavior is explained by the fact that all but the last structure are of type 23S prokaryotic ribosomal RNA, originating in various bacterial organisms [14].

Over the entire bpRNA-1m(90) database, a total of 573 structures had a decrease in numbers of pseudoknots when analyzed with the  $\tau = \infty$  versus  $\tau = 0$  methods. Of these, 531 structures decreased in quantity of pseudoknots by 1, 41 structures by 2, and 1 structure decreased by 3 (bpRNA\_CRW\_55316, *Plasmodium falciparum*). Of the 41 structures which decreased by 2, there were 6 unique RNA types with 36 of them being of type 23S. Structures that changed from having a nonzero quantity of pseudoknots to zero pseudoknots are shown in Table 2. Of these, one structure (*Homo sapiens*) decreased from 2 to 0 pseudoknots. All other structures decreased from 1 to 0 pseudoknots.

**4.1. Maximum Values of  $\tau$  and Persistence of Partitions.** Let  $\tau \geq 1$ . As distances between chords are finite, there is a minimum value of  $\tau$ , say  $\tau_m$ , such that for any  $\tau_* \geq \tau_m$  the  $\tau_*$ -segment partition and the  $\tau_m$ -segment partition are identical. The quantity  $\tau_m$  is precisely the minimum value of  $\tau$  such that the  $\tau$ -segment partition is equivalent to the augmented segment partition. For all structures in bpRNA-1m(90), we calculate  $\tau_m$  by first finding the augmented segment partition, and then finding the  $\tau$ -segment partition for each  $\tau > 0$  until the  $\tau$ -segment partition is equal to the augmented segment partition. We find that the average  $\tau_m$  is 13.035 and the median is 8. The mean absolute deviation is 10.96 and the median absolute deviation is 2. There are 323 structures with  $\tau_m$  at least 17, and 33 structures with  $\tau_m$  at least 100.

Structures with large  $\tau_m$  contain correspondingly large bulges and internal loops; large  $\tau_m$  results from large gaps between chords which are nested but not  $\tau$ -near for many values of  $\tau$ . See for example Figure 11. We verified this by keeping track of  $\tau$ -segment partitions during the process of calculating  $\tau_m$ . In sum, persistent  $\tau$ -segment partitions are indicative of large bulges and internal loops.

**4.2. Classifying Bases.** We implemented the `classifyBases` routine with  $\tau = 0$  to analyze secondary structures and compare quantities obtained from the bpRNA-1m database (see also [12, Figure 8b]). Quantities are shown in Table 12 (left). We observed slight discrepancies in pseudoknot type counts, though the general shape of the distribution is the same. The discrepancies with the  $\tau = 0$  method arise from

the labeling of multiloops and external loops, as some structures in [12] have bases in the external loop that are labeled as a multiloop base. This is a bpRNA software bug that has since been fixed in a fork of [10]. The structure bpRNA\_CRW\_10025 is one example in which the secondary structure labeling is incorrect; bases 357, 385-393, 433-440, 604-633, 690-695, and 728-742 are labeled as part of a multiloop in bpRNA-1m, but by our definition they are part of the exterior loop. Note that the counts of pseudoknot types in Figure 12 (left) differ only when an exterior loop or multiloop is part of the pseudoknot type. Figure 12 (right) shows a comparison of pseudoknot type counts between the  $\tau = 0$  method and the  $\tau = \infty$  method.

**4.3. Calculation of Genus and Clique Numbers.** After implementing the  $\tau$ -segment graph method with  $\tau = 0$  and  $\tau = \infty$  we calculated the genus and maximum clique numbers of  $D_S$  and  $D_A$  and the segment graphs  $G_S$  and  $G_A$  respectively in MATLAB. The results are reported in Figure 13, Table 3, and Figure 14. Out of 3,320 segment graphs, 3,208 are forests, and out of 3,320 augmented segment graphs, 3,210 are forests. Table 4 shows the frequency of forests with a given maximum tree size. This is important to note in the context of using genus and vertex covers for pseudoknot quantification. By Theorem 2.6, if the intersection graph of a chord diagram  $D$  is a forest  $F$ , then the genus  $\gamma(D)$  is equal to  $\beta(F)$ . That is, for acyclic intersection graphs, an increase in genus implies an increase in the vertex cover. From this, we see that the genus of a corresponding chord diagram of an RNA structure is a robust quantifier of pseudoknot complexity. This is further supported by the bubble charts in Figure 14.

#### ACKNOWLEDGEMENTS

The authors were partially supported by the National Science Foundation, DMS-2204148 and The Thomas F. and Kate Miller Jeffress Memorial Trust, Bank of America, Trustee.

R. Ibrahim, DEPARTMENT OF MATHEMATICS, LAFAYETTE COLLEGE, EASTON, PA 18042  
*E-mail address*, R. Ibrahim: [ibrahimr@lafayette.edu](mailto:ibrahimr@lafayette.edu)

A. H. Moore, DEPARTMENT OF MATHEMATICS & APPLIED MATHEMATICS, VIRGINIA COMMONWEALTH UNIVERSITY, RICHMOND, VA 23284  
*E-mail address*, A. H. Moore: [moorea14@vcu.edu](mailto:moorea14@vcu.edu)

## APPENDIX

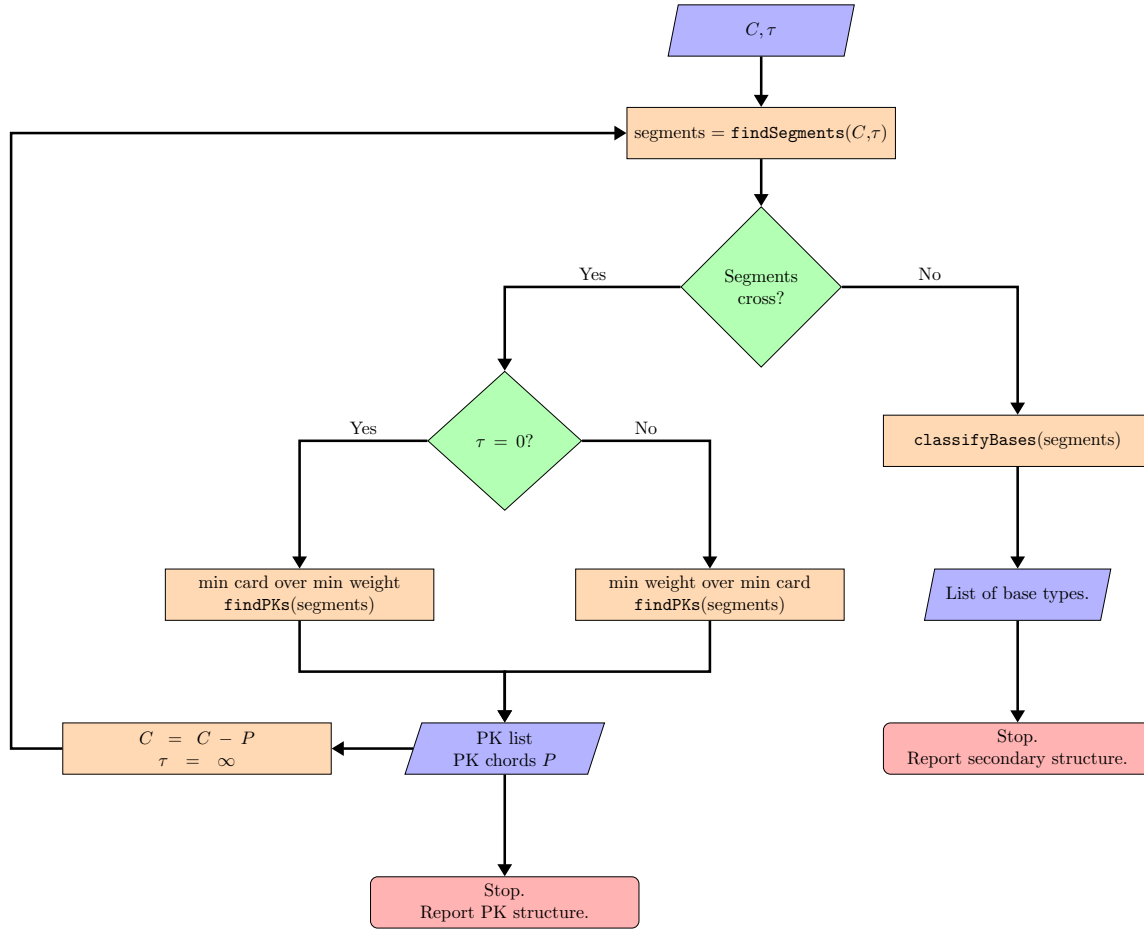


FIGURE 8. A flow chart summarizing the the pseudoknot and secondary structure identification and quantification process described in Section 3.2.

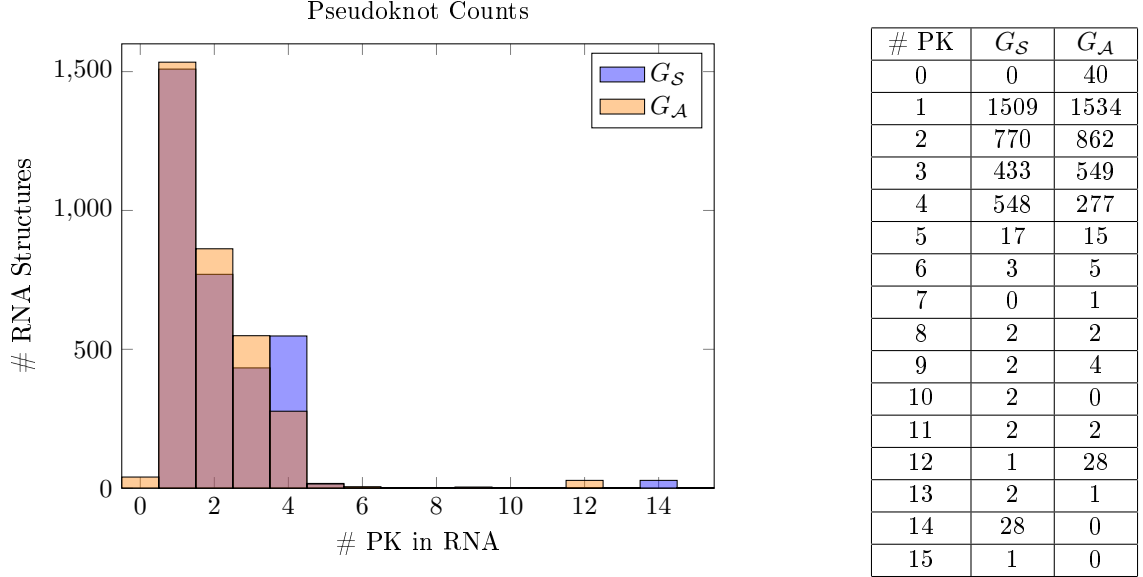


FIGURE 9. A comparison of total number of pseudoknots between the segment graph method ( $\tau = 0$ ) and augmented segment graph method ( $\tau = \infty$ .) For each method, the number of structures containing  $r$  pseudoknots is given.

ID	Domain	Organism	$G_S$	$G_A$	ID	Domain	Organism	$G_S$	$G_A$
CRW_55315	Eukaryota	Euglena gracilis	15	13					
CRW_55268	Bacteria	Acinetobacter calcoaceticus	14	12	CRW_55297	Bacteria	Listeria monocytogenes	14	12
CRW_55269	Bacteria	Aeromonas hydrophila	14	12	CRW_55298	Bacteria	Listeria monocytogenes	14	12
CRW_55271	Bacteria	Bartonella bacilliformis	14	12	CRW_55299	Bacteria	Mycoplasma genitalium	14	12
CRW_55275	Bacteria	Burkholderia mallei	14	12	CRW_55303	Bacteria	Neisseria gonorrhoeae	14	12
CRW_55276	Bacteria	Bordetella pertussis	14	12	CRW_55305	Bacteria	Pseudomonas aeruginosa	14	12
CRW_55279	Bacteria	Clostridium botulinum B	14	12	CRW_55306	Bacteria	Plesiomonas shigelloides	14	12
CRW_55283	Bacteria	Citrobacter freundii	14	12	CRW_55307	Bacteria	Ruminobacter amylophilus	14	12
CRW_55284	Bacteria	Campylobacter jejuni	14	12	CRW_55308	Bacteria	Rickettsia prowazekii (str. Madrid E)	14	12
CRW_55285	Bacteria	Chlamydomonada psittaci 6BC	14	12	CRW_55312	Bacteria	Staphylococcus carnosus	14	12
CRW_55287	Bacteria	Deinococcus radiodurans	14	12	CRW_55313	Bacteria	Thermotoga maritima	14	12
CRW_55290	Bacteria	Enterococcus faecalis	14	12	CRW_55314	Eukaryota	Chlamydomonas reinhardtii	14	12
CRW_55291	Bacteria	Erysipelothrix rhusiopathiae (str. 715)	14	12	CRW_55317	Eukaryota	Spinacia oleracea	14	12
CRW_55292	Bacteria	Haemophilus influenzae (operons A-F)	14	12	CRW_55338	Eukaryota	Cyanophora paradoxa	14	12
CRW_55295	Bacteria	Leptospira interrogans	14	12	CRW_55270	Bacteria	Bacillus anthracis	13	11
CRW_55296	Bacteria	Lactococcus lactis	14	12	PDB_647	Bacteria	Oceanobacillus ihayensis	13	11

TABLE 1. The 31 RNA structures with at least 13 pseudoknots when analyzed with the segment graph method. The rightmost column compares the number of pseudoknots in each structure via the  $\tau$ -segment graph method with  $\tau = \infty$ . All structures with the exception of PDB\_647 are RNA type 23S ribosomal RNA.

ID	Domain	Organism	Length	Method	ID	Domain	Organism	Length	Method
CRW_1213	Bacteria	Actinomyces israelii	734	CSA					
CRW_1219	Bacteria	Actinomyces israelii	1145	CSA	CRW_4401	Bacteria	Streptomyces mobaraensis	1197	CSA
CRW_1563	Bacteria	Clavibacter sp. R1.2_cr	476	CSA	CRW_4409	Bacteria	Streptomyces olivoreticuli	1216	CSA
CRW_1725	Bacteria	Arthrobacter sp.	300	CSA	CRW_4416	Bacteria	Streptomyces salmonis	1136	CSA
CRW_17723	Bacteria	Lachnospira multipara	977	CSA	CRW_4449	Bacteria	coryneform actinomycete B755	679	CSA
CRW_17729	Bacteria	Moorella thermoautotrophica	869	CSA	CRW_4908	Bacteria	Acidocella facilis	922	CSA
CRW_17730	Bacteria	Moorella thermoautotrophica	821	CSA	CRW_4910	Bacteria	Acidiphilium angustum	977	CSA
CRW_17811	Bacteria	Thermoanaerobacter acetothylacus	770	CSA	CRW_4918	Bacteria	Acidiphilium sp.	944	CSA
CRW_17823	Bacteria	Thermoanaerobacter ethanolicus	650	CSA	CRW_7455	Bacteria	unidentified eubacterium 37SW-1	277	CSA
CRW_17834	Bacteria	Thermoanaerobacterium thermosulfurigenes	930	CSA	CRW_7488	Bacteria	Proteobacteria sp	484	CSA
CRW_20267	Bacteria	Marigold phyllody phytoplasma	1015	CSA	CRW_7494	Bacteria	uncultured alpha proteobacterium	410	CSA
CRW_20554	Bacteria	Mycoplasma collis	372	CSA	CRW_7502	Bacteria	uncultured alpha proteobacterium	222	CSA
CRW_20606	Bacteria	Beet leafhopper transmitted virescence phytoplasma	700	CSA	CRW_7614	Bacteria	Nitrobacter sp.	452	CSA
CRW_20626	Bacteria	Potato witches'-broom phytoplasma	658	CSA	CRW_7802	Bacteria	Rhodovulum euryhalinum	1138	CSA
CRW_20629	Bacteria	Paulownia witches'-broom phytoplasma	698	CSA	CRW_7938	Bacteria	Sphingomonas asaccharolytica	629	CSA
CRW_3719	Bacteria	Actinomycetales sp.	472	CSA	CRW_8046	Bacteria	uncultured alpha proteobacterium	751	CSA
CRW_3726	Bacteria	Actinomycetales sp.	503	CSA	CRW_8048	Bacteria	uncultured alpha proteobacterium	730	CSA
CRW_3729	Bacteria	Actinomycetales sp.	502	CSA	CRW_8050	Bacteria	uncultured alpha proteobacterium	690	CSA
CRW_3732	Bacteria	Actinomycetales sp.	478	CSA	PDB_567	artificial sequences	synthetic construct	35	X-RAY
CRW_4109	Bacteria	Mycobacterium xenopi	942	CSA	PDB_512	Eukaryota	Homo sapiens	12	X-RAY
CRW_4363	Bacteria	Streptomyces abikoensis	1177	CSA					

TABLE 2. The 40 RNA structures with nonzero quantity of pseudoknots when analyzed with the segment graph method but zero pseudoknots using  $\tau = \infty$  segment graph method. One structure (*Homo sapiens*) decreased from 2 to 0 pseudoknots. All other structures decreased from 1 to 0 pseudoknots.

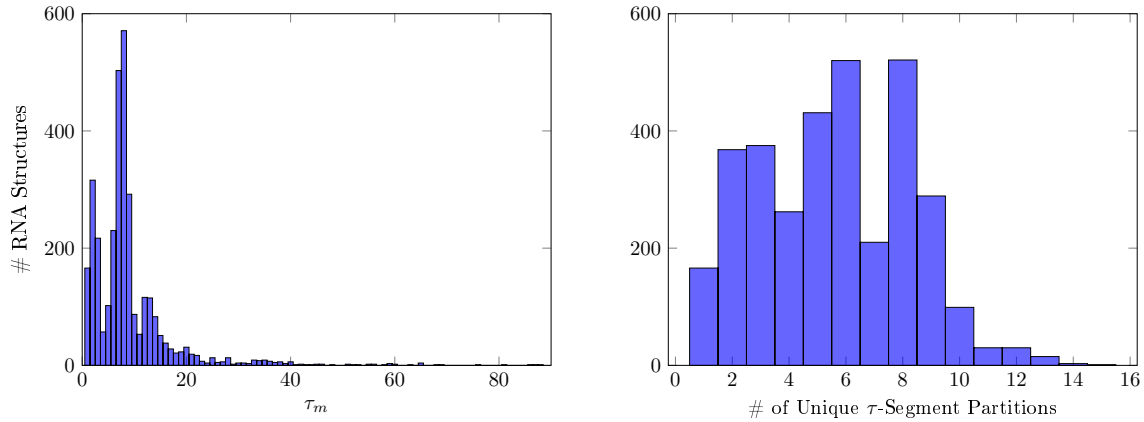


FIGURE 10. Left: Distribution of  $\tau_m$  over all RNA structures in bpRNA-1m(90) restricted to values of  $\tau_m$  within one and a half standard deviations from the mean. Right: Distribution of number of unique  $\tau$ -segment partitions.

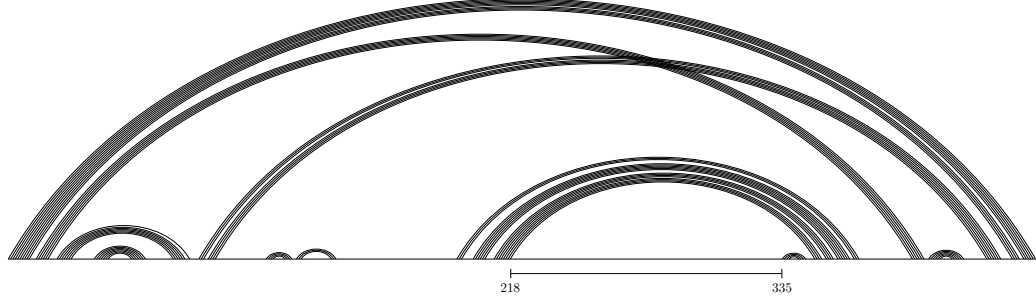
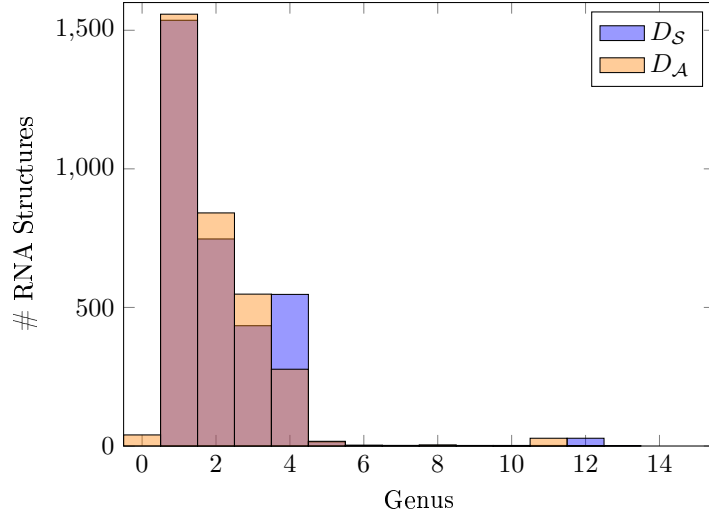


FIGURE 11. The RNA structure bpRNA\_RFAM\_4761 is an example of a structure with large  $\tau_m = 117$  resulting from an interior loop.

PK type	bpRNA-1m	$G_S$
E-E	0	0
E-X	0	0
X-X	0	0
B-B	6	6
I-I	9	9
B-E	10	10
I-M	49	48
E-M	53	53
I-X	57	58
E-I	64	64
M-X	100	104
B-I	153	153
B-X	158	169
H-I	194	194
H-H	261	261
B-M	377	366
E-H	588	588
H-X	670	847
M-M	711	707
B-H	1826	1826
H-M	1878	1701

PK type	$G_S$	$G_A$
E-E	0	0
E-X	0	0
X-X	0	0
B-B	6	5
I-I	9	7
B-E	10	8
I-M	48	48
E-M	53	54
I-X	58	60
E-I	64	57
M-X	104	102
B-I	153	82
B-X	169	5
H-I	194	196
H-H	261	258
B-M	366	9
E-H	588	584
M-M	707	705
H-X	847	845
H-M	1701	1699
B-H	1826	1824

FIGURE 12. (Left) A comparison of counts of pseudoknot types reported in bpRNA-1m versus our  $G_S$  method with  $\tau = 0$ . Discrepancies result from a known bug in the bpRNA software. (Right) A comparison of counts of pseudoknot types using our  $G_S$  and  $G_A$  methods with  $\tau = 0$  and  $\tau = \infty$ , respectively.



Genus	$D_S$	$D_A$
0	0	40
1	1536	1558
2	747	841
3	434	548
4	547	277
5	15	17
6	3	2
7	2	2
8	2	4
9	2	0
10	1	2
11	2	28
12	28	1
13	1	0

FIGURE 13. A comparison of the genus of the segment graph and the augmented segment graph.

$\omega(G)$	$G_S$	$G_A$
1	0	40
2	3214	3177
3	60	96
4	46	7

TABLE 3. The frequency of clique numbers for segment and augmented segment graphs.

Max Tree Order	$G_S$	$G_A$
1	0	40
2	1449	1410
3	998	999
4	637	637
5	27	27
6	71	71

Max Tree Order	$G_S$	$G_A$
7	6	6
8	7	7
9	7	7
10	1	1
11	3	3
12	2	2

TABLE 4. For each segment and augmented segment graph that is a forest, we calculate the maximum order of a tree component.

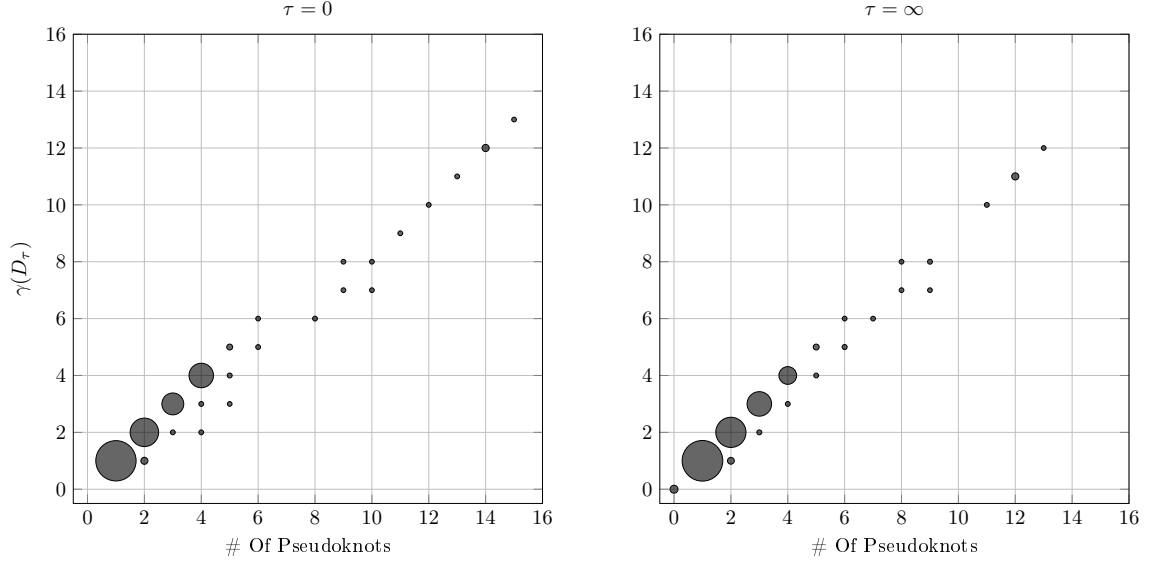


FIGURE 14. Two bubble chart comparisons. Left: Comparison between genera of the chord diagrams  $D_S$  and pseudoknot count via the  $\tau = 0$  method. Right: Comparison between genera of the chord diagrams  $D_A$  and pseudoknot count via the  $\tau = \infty$  method.



## REFERENCES

- [1] Christos Andrikos, Evangelos Makris, Angelos Kolaitis, Georgios Rassias, Christos Pavlatos, and Panayiotis Tsanakas. Knotify: An efficient parallel platform for RNA pseudoknot prediction using syntactic pattern recognition. *Methods and Protocols*, 5(1), 2022.
- [2] Giorgio Benedetti and Stefano Morosetti. A graph-topological approach to recognition of pattern and similarity in RNA secondary structures. *Biophysical Chemistry*, 59(1):179–184, 1996.
- [3] 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, 01 2000.
- [4] Berk Birand. Bron-kerbosch maximal independent set and maximal clique algorithms, 2014.
- [5] Michael Bon, Graziano Vernizzi, Henri Orland, and A. Zee. Topological classification of RNA structures. *Journal of Molecular Biology*, 379(4):900–911, 2008.
- [6] Ian Brierley, Simon Pennell, and Robert J. C. Gilbert. Viral RNA pseudoknots: versatile motifs in gene expression and replication. *Nature Reviews Microbiology*, 5(8):598–610, August 2007.
- [7] Coen Bron and Joep Kerbosch. Algorithm 457: finding all cliques of an undirected graph. *Commun. ACM*, 16(9):575–577, 09 1973.
- [8] J. Brown. The ribonuclease P database. *Nucleic Acids Research*, 26(1):351–352, January 1998.
- [9] Jamie J Cannone, Sankar Subramanian, Murray N Schnare, James R Collett, Lisa M D’Souza, Yushi Du, Brian Feng, Nan Lin, Lakshmi V Madabusi, Kirsten M Müller, Nupur Pande, Zhidi Shang, Nan Yu, and Robin R Gutell. The Comparative RNA Web (CRW) Site: an online database of comparative sequence and structure information for ribosomal, intron, and other RNAs. *BMC Bioinformatics*, 3(1), January 2002.
- [10] Padideh Danaee. bpRNA. <https://github.com/padidehdanaee/bpRNA>, 2018.
- [11] Padideh Danaee, Mason Rouches, Michelle Wiley, Dezhong Deng, Liang Huang, and David Hendrix. bpRNA-1m: A database of single-molecule RNA secondary structures annotated by bpRNA. <https://bprna.cgrb.oregonstate.edu/about.php>, 2018.
- [12] Padideh Danaee, Mason Rouches, Michelle Wiley, Dezhong Deng, Liang Huang, and David Hendrix. bpRNA: large-scale automated annotation and analysis of RNA secondary structure. *Nucleic Acids Research*, 46(11):5381–5394, May 2018.
- [13] P. Dumas, D. Moras, C. Florentz, R. Giegé, P. Verlaan, A. Van Belkum, and C. W.A. Pleij. 3-D Graphics Modelling of the tRNA-Like 3’-End of Turnip Yellow Mosaic Virus RNA: Structural and Functional Implications. *Journal of Biomolecular Structure and Dynamics*, 4(5):707–728, April 1987.
- [14] David Elliott and Michael Lodomery. *Molecular Biology of RNA*. Oxford University Press, Oxford, 2023.
- [15] S. Griffiths-Jones. Rfam: an RNA family database. *Nucleic Acids Research*, 31(1):439–441, January 2003.
- [16] Rayan K. Ibrahim. bprnray. <https://github.com/raymaths/bpRNRay>, 2025.
- [17] F. Juhling, M. Morl, R. K. Hartmann, M. Sprinzl, P. F. Stadler, and J. Putz. tRNAdb 2009: compilation of tRNA sequences and tRNA genes. *Nucleic Acids Research*, 37(Database):D159–D162, January 2009.
- [18] Marcel Kucharík, Ivo L. Hofacker, Peter F. Stadler, and Jing Qin. Pseudoknots in RNA folding landscapes. *Bioinformatics*, 32(2):187–194, October 2015.

- [19] Terry A. McKee and F. R. McMorris. *Topics in intersection graph theory*. SIAM Monographs on Discrete Mathematics and Applications. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1999.
- [20] Gadi Moran. Chords in a circle and linear algebra over  $\text{GF}(2)$ . *J. Combin. Theory Ser. A*, 37(3):239–247, 1984.
- [21] R. C. Penner and Michael S. Waterman. Spaces of RNA secondary structures. *Adv. Math.*, 101(1):31–49, 1993.
- [22] Christian Reidys. *Combinatorial computational biology of RNA*. Springer, New York, New York, NY, 2011. Pseudoknots and neutral networks.
- [23] Christian M. Reidys, Fenix W. D. Huang, Jørgen E. Andersen, Robert C. Penner, Peter F. Stadler, and Markus E. Nebel. Topology and prediction of RNA pseudoknots. *Bioinformatics*, 27(8):1076–1085, 02 2011.
- [24] K. Rietveld, R. Van Poelgeest, C.W.A. Pleij, J.H. Van Boom, and L. Bosch. The tRNA-Uke structure at the 3' terminus of turnip yellow mosaic virus RNA. Differences and similarities with canonical tRNA. *Nucleic Acids Research*, 10(6):1929–1946, 1982.
- [25] M. A. Rosenblad. Srpdb: Signal recognition particle database. *Nucleic Acids Research*, 31(1):363–364, January 2003.
- [26] William R. Schmitt and Michael S. Waterman. Linear trees and RNA secondary structure. *Discrete Appl. Math.*, 51(3):317–323, 1994.
- [27] Alexander Serganov. Crystal structure of the thermotoga maritima lysine riboswitch bound to s-(2-aminoethyl)-l-cysteine, 2008.
- [28] Alexander Serganov, Lili Huang, and Dinshaw J. Patel. Structural insights into amino acid binding and gene control by a lysine riboswitch. *Nature*, 455(7217):1263–1267, September 2008.
- [29] Wenjie Shu, Xiaochen Bo, Zhiqiang Zheng, and Shengqi Wang. A novel representation of RNA secondary structure based on element-contact graphs. *BMC Bioinformatics*, 9(1), April 2008.
- [30] David W Staple and Samuel E Butcher. Pseudoknots: RNA structures with diverse functions. *PLoS Biology*, 3(6):e213, June 2005.
- [31] Gary M. Studnicka, Georgia M. Rahn, Ian W. Cummings, and Winston A. Salser. Computer method for predicting the secondary structure of single-stranded rna. *Nucleic Acids Research*, 5(9):3365–3388, 1978.
- [32] Ignacio Tinoco, Olke C. Uhlenbeck, and Mark D. Levine. Estimation of secondary structure in ribonucleic acids. *Nature*, 230(5293):362–367, April 1971.
- [33] Michael S. Waterman. Secondary structure of single-stranded nucleic acids. In *Studies in foundations and combinatorics*, volume 1 of *Adv. Math. Suppl. Stud.*, pages 167–212. Academic Press, New York-London, 1978.
- [34] Michael S. Waterman. Combinatorics of RNA hairpins and cloverleaves. *Stud. Appl. Math.*, 60(2):91–96, 1979.
- [35] M.S. Waterman and T.F. Smith. RNA secondary structure: a complete mathematical analysis. *Mathematical Biosciences*, 42(3):257–266, 1978.
- [36] J.D. Watson, T.A. Baker, S.P. Bell, A. Gann, M. Levine, and R. Losick. *Molecular Biology of the Gene*. Pearson Education, San Francisco, 2013.
- [37] C. Zwieb. tmRDB (tmRNA database). *Nucleic Acids Research*, 31(1):446–447, January 2003.