



Predicting Large RNA-Like Topologies by a Knowledge-Based Clustering Approach

Naoto Baba^{1,2}, Shereef Elmetwaly¹, Namhee Kim¹ and Tamar Schlick^{1,3}

1 - Department of Chemistry and Courant Institute of Mathematical Sciences, New York University, 251 Mercer Street, New York, NY 10012, USA

2 - Department of Chemistry, Nagoya University, Furo-cho, Chikusa-ku, Nagoya, Aichi 464-8601, Japan

3 - NYU-ECNU Center for Computational Chemistry at NYU Shanghai, 3663 Zhongshan Road North, Shanghai, 200062, China

Correspondence to Tamar Schlick: schlick@nyu.edu

<http://dx.doi.org/10.1016/j.jmb.2015.10.009>

Edited by A. Pyle

Abstract

An analysis and expansion of our resource for classifying, predicting, and designing RNA structures, RAG (RNA-As-Graphs), is presented, with the goal of understanding features of RNA-like and non-RNA-like motifs and exploiting this information for RNA design. RAG was first reported in 2004 for cataloging RNA secondary structure motifs using graph representations. In 2011, the RAG resource was updated with the increased availability of RNA structures and was improved by utilities for analyzing RNA structures, including substructuring and search tools. We also classified RNA structures as graphs up to 10 vertices (~200 nucleotides) into three classes: existing, RNA-like, and non-RNA-like using clustering approaches. Here, we focus on the tree graphs and evaluate the newly founded RNAs since 2011, which also support our refined predictions of RNA-like motifs. We expand the RAG resource for large tree graphs up to 13 vertices (~260 nucleotides), thereby cataloging more than 10 times as many secondary structures. We apply clustering algorithms based on features of RNA secondary structures translated from known tertiary structures to suggest which hypothetical large RNA motifs can be considered “RNA-like”. The results by the PAM (Partitioning Around Medoids) approach, in particular, reveal good accuracy, with small error for the largest cases. The RAG update here up to 13 vertices offers a useful graph-based tool for exploring RNA motifs and suggesting large RNA motifs for design.

© 2015 Elsevier Ltd. All rights reserved.

Introduction

It is now well appreciated that RNA molecules have essential roles in the regulation of gene expression and signal recognition [1–4] besides their widely known roles in protein synthesis by mRNA, tRNA, and rRNA. The functionalities of RNAs are made possible by large variations of secondary and tertiary motifs. Unlike proteins, where structural genomics initiatives have been advancing for decades [5,6], systematic connections between RNA structures and their biological roles remain largely unclear. Thus, improvements in the connection between RNA's structure and its functionality can help advance our understanding of RNAs and the design of new RNAs.

The secondary structure of RNA, less complex than its tertiary structure, is already a good starting point for a structural/functional analysis. Secondary

structures, in particular, are amenable to mathematical analysis by graph theory. Graph theory is a well-established field of mathematics, which has been used extensively in a variety of economic, social, engineering, biological, and medical contexts to describe and analyze complex networks [7–10]. Shareability networks have been used recently, for example, to analyze cab sharing in New York City and to propose a 40% reduction in traffic and pollution due to simple sharing of cabs [11]. We utilize graph theory here to analyze RNA secondary structures: we transform RNA secondary structures into graph vertices and edges to express RNAs as coarse-grained objects, thereby forgoing a detailed atomic-level representation (see Fig. 1). Applying graph theory to compare the 2D (2-dimensional) graphical representations has already shown to be useful in some projects [12–14].

In 2004, we developed and launched the RAG (RNA-As-Graphs) Web resource[†]. This framework catalogs all possible RNA 2D topologies up to 10 vertices and classifies them as existing or hypothetical, with the latter divided into RNA-like (“non-existing but RNA-like”) and non-RNA-like (“non-existing and not RNA-like”) [15], by clustering features of RNA secondary structures as tree and dual graphs by means of graph theory. The graphical information extracted is in the form of the adjacency and Laplacian matrices, which describe graph connections, and the clustering is performed by their vertex number and eigenvalue spectrum (see [Materials and Methods](#)).

The many applications of RAG, as reviewed recently [16–18], include the prediction of RNA-like topologies [19–22], prediction of non-coding RNA [23,24], computational modeling of the *in vitro* selection process for RNA design [25–27], analysis of large viral RNA [28,29], analysis and design of riboswitches [30,31], graph partitioning to explore RNA modularity [16,17,32], and prediction of 3D (3-dimensional) RNA topologies [33,34].

Many new RNA databases have been developed since 2004. For example, the RNA family database Rfam [35] displays consensus secondary structures for 1372 families of RNA [36], and the RNA STRAND database catalogs 4666 secondary structures determined by comparative sequence analysis, NMR data, and X-ray crystallography [37]. This growth allowed us to extend RAG and propose an improved classification in 2011. In addition, we implemented various improvements to the RAG Web resource such as expanded search tools and a user-friendly interface.

The 2011 update was still limited to tree graphs up to 10 vertices corresponding to about 200 nucleotides of RNA sequences.

In this work, we upgrade the RAG database with new prediction results for RNA-like topologies for large tree graphs up to 13 vertices (~260 nucleotides) in length, using an auxiliary graph computation program named *nauty* and *Traces* [38]. This makes RAG’s coverage more than 10-fold greater. We then catalog new existing RNAs from the Protein Data Bank (PDB) database, as of August 2014, for all secondary structures translated from solved experimental structures. Finally, a new prediction for RNA-like motifs is described based on the PAM (Partitioning Around Medoids) clustering approach [39].

Our main achievements consist of the two parts: high accuracy of predicted RNA-like features for the newly found RNAs and our extended RAG for larger topologies based on the current dataset. In [Discussion](#), we elaborate upon the significance of those findings, and we mention the future prospects of clustering for RNAs.

Results

Association of secondary structures to new RNAs

The process of converting an RNA 2D full topology into a tree graph, which is described in the [Materials and Methods](#) section, is automated in RAG [19]. This allowed us to exhaustively inspect the current RNA structures and assign a secondary graph motif to each. Taking RNA structures from the PDB yielded [Fig. 2](#). Many new topologies were identified, even

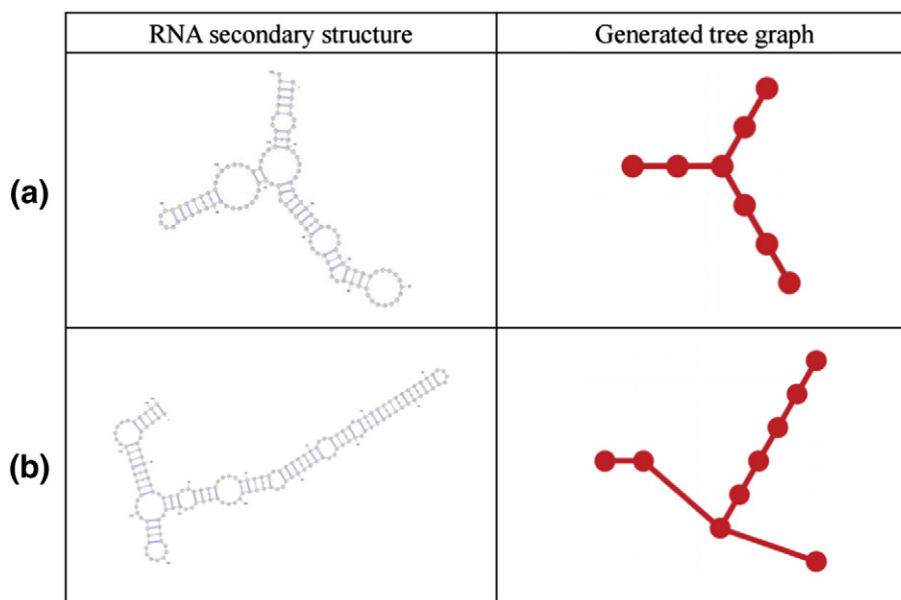

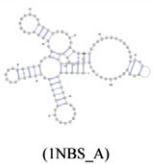

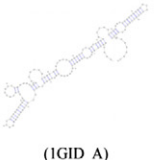
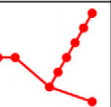

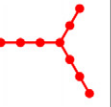


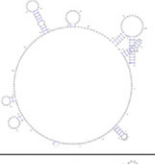
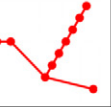
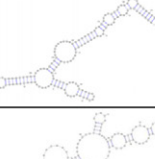
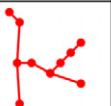
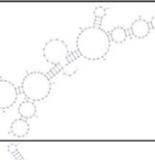

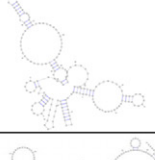

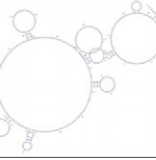


Fig. 1. Conversion from a secondary RNA structure into a planar tree graph. (a) 16S ribosomal RNA (PDB ID 3 J12, chain A) with its tree graph. (b) 80S ribosomal RNA (PDB ID 3IZD, chain A) with its tree graph.

(a)

| Graph ID | RAG motif | RNA 2D Structure | RNA (PDB ID) |
|----------|--|--|--|
| 8_15 |  |  | Ribonuclease P RNA (1NBS_A), 18S ribosomal RNA (3J16_K) |
| 9_2 |  |  | Group I Intron (1GID_A, 1GID_B, 1HR2_B) |
| 9_4 |  |  | 80S ribosomal RNA (3IZD_A) |
| 9_19 |  |  | Signal Recognition Particle (1L9A_B, 1MFQ_A, 2G05_A, 2J37_A) |

(b)

| | | | |
|-------|---|---|--|
| 9_46 |  |  | Ribonuclease P Bacterial A-type (2A2E_A) |
| 10_4 |  |  | M-Box Riboswitch Aptamer Domain (2QBZ_X) |
| 10_19 |  |  | Glycine Riboswitch (3P49_A) |
| 10_45 |  |  | Adenosylcobalami n Riboswitch (4GMA_Z) |
| 10_69 |  |  | Transfer-messenger RNA (tmRNA) (31YR_A) |

from the RNAs that had been identified before our last work, because our current procedure for excision of pseudoknots and separation of multiple chains allows the conversion of the RNA structures that could not be handled previously as tree graphs.

Clustering procedure and current assessment

Early in our RAG project, the two clustering methods, PAM [39] and k -nearest neighbor (k -NN) [40,41], were used for predicting novel RNA topologies based on clustering. Because k -NN considers randomized data for its prediction, we consider it now to be less reliable than PAM.

Indeed, by the procedure described in the [Materials and Methods](#) section below ([Clustering and validation procedure](#)), we obtain 77.27% accuracy from PAM ([Fig. 3](#) and [Table 1](#)) compared to poorer results by k -NN (see Supplementary Material).

Clustering and validation procedure

Overall, our goal is to predict which of the hypothetical tree graphs are RNA-like. To do so, we cluster the data points generated from the tree graphs are clustered into two categories: RNA-like and non-RNA-like. Two very different clustering approaches can be considered: k -nearest neighbor (k -NN) [40,41] and partitioning around medoids (PAM) [39]. The former uses training data while the latter does not.

The k -NN algorithm classifies a point based on k closest training data points: a point is classified by a majority vote of its neighbors, with the point being assigned to the class most common among its k nearest neighbors [40,41]. However, due to the lack of existing motifs for higher vertices, we use all existing motifs and the same number of randomly selected non-existing motifs as a training set. Because of this randomness, we employed 10 trials by varying the set of random non-existing data.

Once a training set is given, cross-validation is one of several approaches for estimating how well the model might perform on future data. One effective cross-validation method is called leave-one-out cross-validation (LOOCV) [42]. As its name suggests, LOOCV leave-one-out cross-validation leaves one data item from the training set and performs a clustering to this single isolated data point by the training set which that now lacks that item. This process is repeated for each data item, and the reliability of the prediction is measured by comparison to confirmed RNA-like and non-RNA-like motifs.

PAM, on the other hand, requires no training set. PAM partitions all data (existing and hypothetical

Fig. 2. List of newly found motifs and their associated secondary structures of RAG graphs. For up through 10-vertex graphs, nine new motifs have been found since our last update.

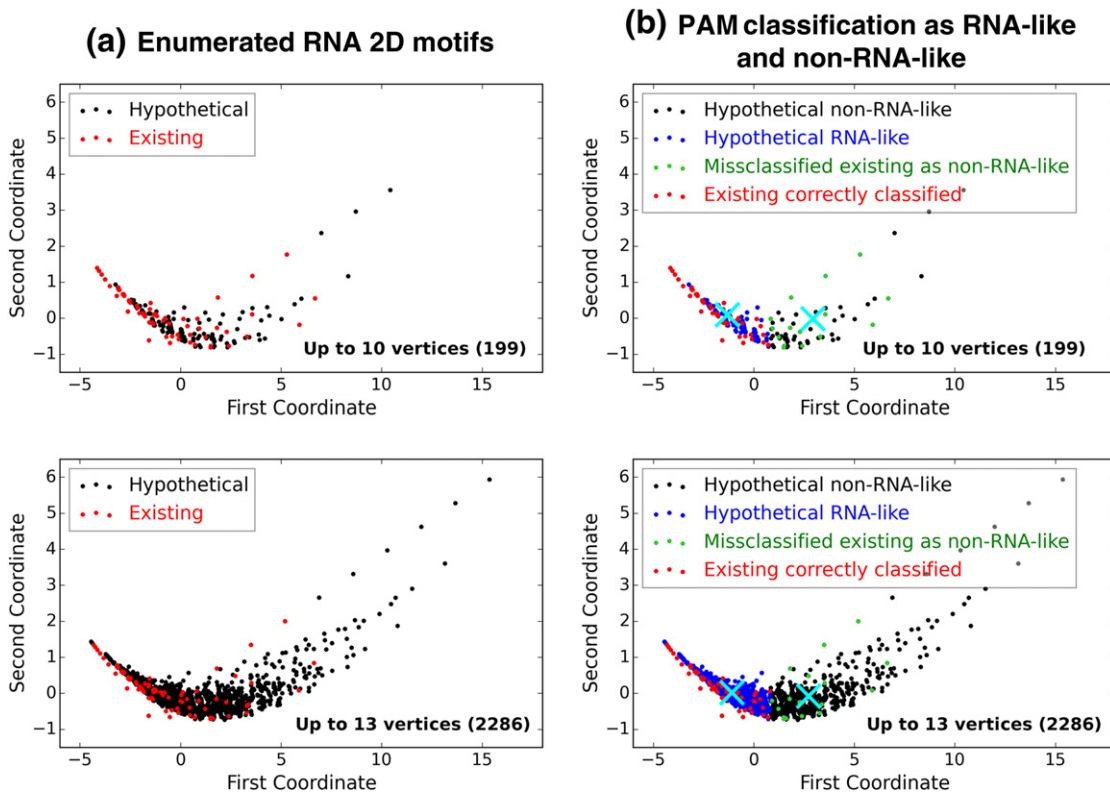


Fig. 3. Plot of PAM clustering result. (a) Enumerated RNA 2D motifs up to 10 vertices (upper) and 13 vertices (lower): the x-axis and the y-axis are the variables reduced by the MDS as described in Deduction of characteristic information from the Laplacian spectra. Red indicates existing RNAs. (b) PAM classification as RNA-like and non-RNA-like up to 10 vertices (upper) and 13 vertices (lower): the two medoids, or centers, of PAM are indicated by X. Most existing RNAs (65 of 84 existing RNAs) are confirmed as the RNA-like group (red) but 19 are classified as non-RNA-like (green). Hypothetical RNAs are further divided and predicted into RNA-like (blue) and non-RNA-like (black) by the PAM clustering approach.

graph features) in an “*ab initio*” manner to predict two groups (RNA-like and non-RNA-like) that are maximally separated [39]. Thus, PAM clusters the data into these two groups, each with its center or medoid, by minimizing the distances within groups and maximizing the distance between groups.

The fact that the PAM requires no training set makes the validation fairly straightforward. We simply perform PAM clustering on the current dataset and calculate the accuracy naturally by

$$\frac{(\text{Total number of existing RNAs predicted correctly as RNA-like})}{(\text{Number of known existing RNAs})}$$

We further check and confirm actual existing RNAs predicted as either RNA-like or non-RNA-like graphs (i.e., that we get not just the right number but the right graphs).

High accuracy of RAG prediction on the newly found RNAs

The PAM clustering method classifies the motifs associated with the newly found RNAs in Table 2, as shown in Fig. 3. Many of the newly found RNAs were

categorized as RNA-like by the RAG clustering strategy. Notably, although three motifs were misclassified as non-RNA-like, they all have only one existing RNA; the motifs that have multiple existing RNAs were all correctly classified as RNA-like.

The RNAs that are misclassified are the following: RNA component of bacterial ribonuclease P (PDB ID 2A2E, chain A) [43], adenosylcobalamin riboswitch (PDB ID 4GMA, chain Z) [44], and tmRNA-SmpB ribonucleoprotein complex (PDB ID 3IYR, chain A) [45].

Drastically extended RAG for larger topologies and its accuracy based on the current dataset

The number of vertices for RNAs is not limited to 10 because nauty and Traces can generate secondary graphs with more vertices. By integrating this software with our program, we exhaustively created all tree graphs through 13 vertices, which allows the enumeration of much larger sets of topological descriptors. Thus, RAG has extended its coverage by more than 10-fold; RAG in 2011 cataloged 199 secondary graph motifs, but now, the count is 2286, with 2087 graph motifs added. Since the graph motifs with varying

Table 1. Statistics from PAM.

| Vertex | Known | | Predicted | | | | Total |
|--------|----------|--------------|-----------------------------|-------------------------------------|--------------|--------------|-------|
| | Existing | Hypothetical | Existing | | Hypothetical | | |
| | | | RNA-like (correct class) | Non-RNA-like (misclassification) | RNA-like | Non-RNA-like | |
| 3 | 1 | 0 | 1 | 0 | 0 | 0 | 1 |
| 4 | 2 | 0 | 2 | 0 | 0 | 0 | 2 |
| 5 | 3 | 0 | 2 | 1 | 0 | 0 | 3 |
| 6 | 6 | 0 | 4 | 2 | 0 | 0 | 6 |
| 7 | 9 | 2 | 6 | 3 | 2 | 0 | 11 |
| 8 | 16 | 7 | 13 | 3 | 4 | 3 | 23 |
| 9 | 15 | 32 | 12 | 3 | 21 | 11 | 47 |
| 10 | 14 | 92 | 11 | 3 | 60 | 32 | 106 |
| 11 | 8 | 227 | 5 | 3 | 156 | 71 | 235 |
| 12 | 4 | 547 | 4 | 0 | 391 | 156 | 551 |
| 13 | 6 | 1295 | 5 | 1 | 934 | 361 | 1301 |
| Total | 84 | 2202 | 65 | 19 | 1568 | 634 | 2286 |

Existing and hypothetical RNA tree motifs, each divided into RNA-like and non-RNA-like by the PAM clustering approach (see Fig. 3 and Fig. 4). For the 2286 motifs up to 13 vertices, 65 are predicted correctly as RNA-like but 19 are false positives.

numbers of nodes are clustered together in RAG, we can make RNA-like predictions for larger topologies regardless of the lack of larger existing motifs. Such predictions can be evaluated based on the RNAs archived from the PDB, which includes new RNAs in addition to the others that we could not represent in

2011. The result is shown as Table 1. The result for 11 vertices is somewhat poor, but there is only one misclassified structure for 13 vertices, and there is no error for 12 vertices. There is only one graph, RAG ID 11_24, with multiple existing RNAs, and it is predicted properly as RNA-like. Table 1 also shows the statistics

Table 2. Newly found RNA motifs and their prediction classes.

| Graph ID | Label | RNA (PDB ID) |
|----------|--------------|---|
| 8_15 | RNA-like | Ribonuclease P RNA (1NBS_A), 18S ribosomal RNA (3J16_K) |
| 9_2 | RNA-like | Group I intron (1GID_A, 1GID_B, 1HR2_B) |
| 9_4 | RNA-like | 80S ribosomal RNA (3IZD_A) |
| 9_19 | RNA-like | Signal recognition particle (1L9A_B, 1MFQ_A, 2GO5_A, 2J37_A) |
| 9_46 | non-RNA-like | Ribonuclease P bacterial A-type (2A2E_A) |
| 10_4 | RNA-like | M-box riboswitch aptamer domain (2QBZ_X) |
| 10_19 | RNA-like | Glycine riboswitch (3P49_A) |
| 10_45 | non-RNA-like | Adenosylcobalamin riboswitch (4GMA_Z) |
| 11_1 | RNA-like | 23S ribosomal RNA (3J5S_A) |
| 11_24 | RNA-like | M-box riboswitch (3PDR_A, 3PDR_X) |
| 11_56 | RNA-like | Ribonuclease P (1U9S_A) |
| 11_89 | non-RNA-like | Transfer-messenger RNA (3IYQ_A) |
| 11_138 | RNA-like | Group 1 intron (3BO4_B) |
| 11_177 | RNA-like | Ribonuclease P (1NBS_B) |
| 11_207 | non-RNA-like | RNase P (3DHS_A) |
| 11_216 | non-RNA-like | Group I intron with a tyrosyl-tRNA synthase (2RKJ_C) |
| 12_150 | RNA-like | Tetrahymena ribozyme (1GRZ_A) |
| 12_286 | RNA-like | 80S ribosomal RNA (3ZEX_E) |
| 12_387 | RNA-like | Group I intron (3IIN_B) |
| 12_392 | RNA-like | Group I intron (3BO2_BCDE) |
| 13_140 | RNA-like | Adenosylcobalamin riboswitch (4GXY_A) |
| 13_181 | RNA-like | Tetrahymena ribozyme (1GRZ_B) |
| 13_1021 | RNA-like | Group I intron (1U6B_CDB) |
| 13_1047 | RNA-like | Group I intron (3BO3_CDB) |
| 13_1154 | non-RNA-like | Group I intron-product complex (1Y0Q_A) |
| 13_1213 | RNA-like | 28S ribosomal RNA (3J16_J) |

For motifs less than or equal to 10 vertices, motifs include updates since our 2011 RAG version. For motifs larger than 10 vertices, motifs are new. Many of the newly found graph motifs are classified as RNA-like. A few of them are misclassified as non-RNA-like, but those motifs only have a single RNA each. For example, there are four RNAs found for ID 9_4, which are RNA-like, but only one for ID 9_46, which is non-RNA-like. The larger RNA motifs more than 11 vertices include only new data. Although there are some misclassified data for 11 vertices, the other results for 12 and 13 nodes are very good. Only one RNA graph, 11_24, has two RNAs, and it is properly predicted as RNA-like.

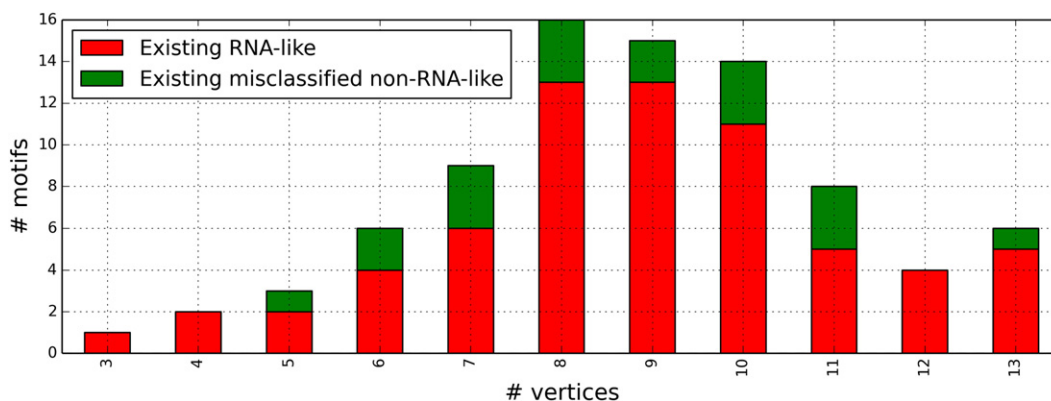


Fig. 4. Numbers of existing RNA-like and existing non-RNA-like vs versus number of vertices. This pictorial view of the statistics obtained in Table 1 and Table 2 reveals that there are more existing RNA-like (properly predicted) topologies than existing misclassified non-RNA-like (incorrectly predicted) topologies for every number of vertices.

for higher vertices, and Fig. 4 visualizes the counts of existing RNA-like and existing misclassified non-RNA-like in Table 1.

Finally, a complete catalog of our RAG data is available. Because of space limitations, only a subset is shown in Fig. 5 for 10-vertex graphs. The full catalog can be found in the Supplementary Material and on our RAG Web site[†].

Discussion

We have extensively updated our RAG database based on the newly discovered RNA structures by exhaustive enumeration of RAG motifs represented as tree graphs up through 13 vertices. Our clustering results show two significant gains: the RAG clustering strategy yields near 80% accuracy for predicting existing RNA topologies, and no motif with multiple existing RNA structures is misclassified. Thus, estimating features of RNA-like structures according to their topological representation may be attractive for RNA design. The predicted RNA-like candidates are good design candidates, as already suggested [15,16,19].

In our previous work [15], we used a build-up approach to predict and identify sequences that fold onto 10 candidate dual graph motifs. Among those 10 candidate motifs, five have since been experimentally determined [16,19]. To design RNA sequences that fold onto the targeted RNA-like topologies, we have used graph partitioning algorithms based on Laplacian eigenvectors [32]. We recently suggested a gap cut approach that partitions a graph into two graphs by the largest gap of the sorted second Laplacian eigenvector μ_2 ; we have illustrated how to use this gap cut partitioning to describe basic modules of RNAs and propose their hierarchical assembly [32].

Figure 6 sketches a design application for RNA-like graphs. Here, we aim to design a large RNA-like graph, RAG ID 11_205. The gap cut suggests

partitioning the graph 11_205 into two substructures, an existing 5_3 corresponding to tRNA (PDB ID 2DU3) and an RNA-like 7_4 graph. The latter graph is further partitioned into two identical existing graphs 4_2 corresponding to the hammerhead ribozyme (PDB ID 1RMN). The assembly of these existing sequences provides a starting candidate sequence for the large RNA corresponding to the target RNA-like graph 11_205. Of course, computational refinements by 2D structure prediction programs combined with thermodynamic and experimental validation are needed for confirmation. However, this systematic design protocol for novel RNA-like topologies could help expand the structural and functional repertoire of RNAs.

Although the RAG classification and prediction described here exhibited good accuracy for predicting existing RNA topologies, many improvements can be envisioned. In addition to eigenvalues, Laplacian eigenvectors could also be useful for graph descriptors. The second eigenvector was shown to be useful for graph partitioning for the discovery of RNA modularity [32]. This kind of approach reveals a connection between RNAs' higher-order structures and their properties. A challenge for the future is to integrate other descriptors and other methods with the current strategy to improve the results.

Conclusion

Focusing on tree graphs, we show our refined RAG classification method to predict well RNA-like and non-RNA-like topologies of secondary structures with near 80% accuracy. We have also expanded the database significantly to larger topologies, adding 10 times as many topologies since the last update. Our analysis suggests that a topology prediction approach can be productive and that it reinforces the idea that the properties of RNAs can be analyzed to a first approximation by means of their secondary structures.

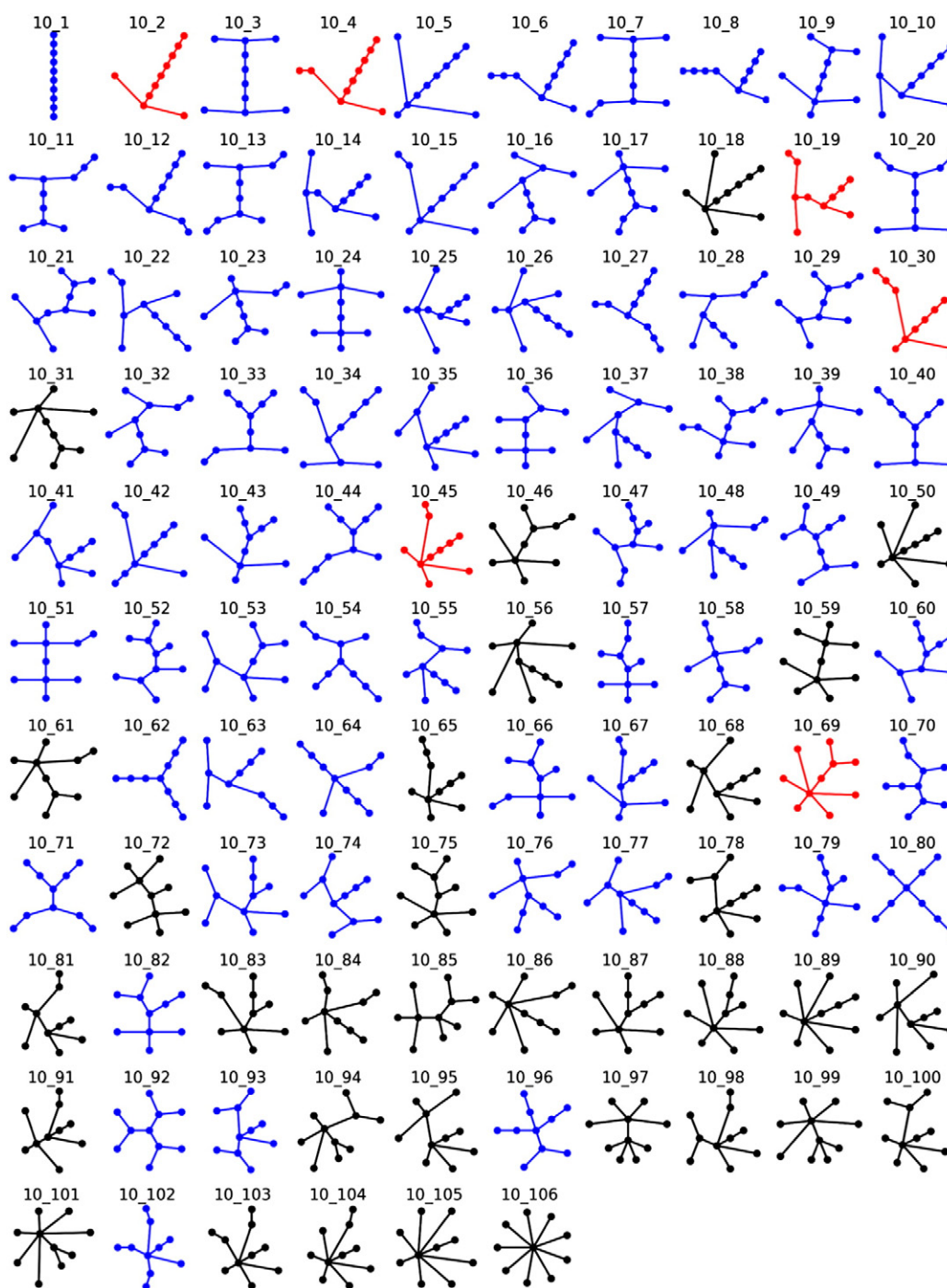


Fig. 5. Illustrative subset of the RAG catalogue. We classify all enumerated graph motifs as existing, RNA-like and non-RNA-like motifs. Existing motifs are colored in red, RNA-like are in blue, and non-RNA-like are in black. The complete version is available in Supplemental Supplementary Material or at <http://www.biomath.nyu.edu/rag/home>.

Materials and Methods

RNA secondary structure data

In our previous works, we used several RNA secondary structure repositories: Rfam [46], PseudoBase++ [47], RNA

STRAND [48], PDB [49], and Nucleic Acid Database [50,51], for cataloging secondary structures that are either fully or partially evaluated by experiment. Here, to analyze the accuracy and efficiency of our RAG clustering strategy for predicting RNA-like motifs, we exclusively collected RNA secondary structures from PDB with untangling of multiple chains so that the structures we classify are all experimentally

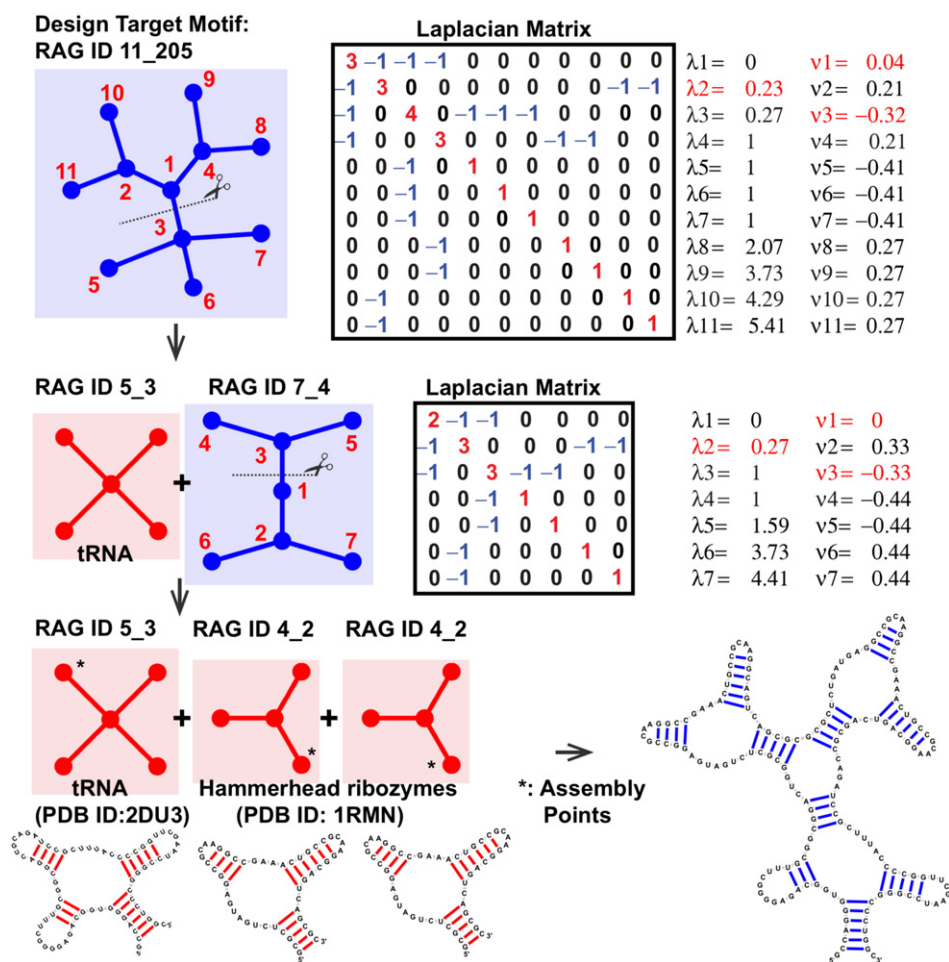


Fig. 6. Design application for RNA-like topologies (example target: RAG ID 11_205). The design procedures using graph partitioning and build-up approaches are shown. In the first row, graph 11_205 (with random vertex numbering), corresponding Laplacian matrix, eigenvalues (λ_2 in red), and the second eigenvector (μ_2) are shown. The largest gap of the sorted elements of μ_2 (vertices 1 and 3) is marked in red. In the second row, two subgraphs (existing graph 5_3 and RNA-like graph 7_4) and gap cut analysis of RNA-like graph 7_4 are shown. The third row shows the assembly procedure: the build-up of three existing modules at the assembly points suggested by gap partitioning produce a candidate RNA with the targeted graph 11_205.

validated. We also include pseudoknot structures, which are translated into non-pseudoknot structures for a representation as tree graphs by removal of extra base pairings composing the pseudoknots. Note that dual graphs, as we have described separately (see Ref. [15] and L. Petingi, N. Kim, and T. Schlick, unpublished results), can be used to model pseudoknotted RNA fully. A simple modification of tree graphs to model pseudoknots was also recently presented and applied for prediction of tertiary structures [18].

RNA tree graph representation

The conversion process from detailed RNA secondary structures into tree graph representations was detailed in our previous works [15,19]. Briefly, RAG considers nucleotide bulges, hairpin loops, internal loops, junctions, and the 3' and 5' ends as vertices, as well as RNA stems as edges (see Fig. 1).

Enumeration of RNA graphs

To classify all existing graph motifs including the experimentally found and those not yet solved experimentally, we generate all possible tree graphs with a given number of vertices. Graph theory offers enumeration methods for describing all possible graphs [52]. Previously, we had used the counting polynomial of Harary–Prins and the figures of graph theory [52], but this scheme for tree graphs was manual; the polynomial gives the number of the graphs but no information about the shape, or topology, of the graphs.

An alternative is the integration of nauty and Traces [38], two programs focused on canonical labeling and automorphism group computations. These programs can exhaustively produce all desired tree graphs. The completeness of the graph generation is verified by two requirements: the number of generated graphs should match the result of the counting polynomial of Harary–Prins, and there should be no isomorphic graphs, which is confirmed by NetworkX [53].

Thus, we ensure that all the non-isomorphic graphs are generated. This effective combination allows us to extend RAG significantly by adding 235, 551, and 1301 tree graphs for 11, 12, and 13 vertices, respectively.

Topological descriptors of RNA graphs: Laplacian spectra

To order all the graphs by their features, we use the second eigenvalue λ_2 of the Laplacian matrix, a matrix that describes graph connections. The other eigenvalues are associated with a spectral decomposition associated with the graph, useful for many applications, for example, graph partitioning by the second eigenvector [32].

To define the Laplacian matrix, we define the $n \times n$ adjacency matrix for an n -node graph where the non-diagonal entries a_{ij} are 1 if there is an edge between vertex i and j and are 0 otherwise.

The Laplacian matrix (L) is defined by $L = D - A$, where D is the diagonal matrix whose diagonal elements a_{ii} specify the degree of connectivity of vertex i . Thus, for example, a straight-line-shaped graph with three vertices has graph ID 3_1 in the RAG terminology and corresponding D , A , and L matrices as follows:

$$D = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{pmatrix}, A = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix},$$

$$L = D - A = \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{pmatrix}$$

Note that the spectrum of the Laplacian matrix is independent of the labeling of graph vertices because a change in labeling can be accomplished by the elementary operations on the matrices and the elementary operations do not alter their eigenvalues. Thus, if the spectra of the Laplacian matrices of two graphs are different, the graphs are also different. Although identical spectra can be associated with different graph topologies, this situation is rare [38].

The pattern of a graph's connectivity is related to its eigenvalue spectrum (spectral graph theory) [54]. The second smallest eigenvalue, λ_2 , for example, is called the algebraic connectivity and measures the graph's compactness: a linear chain has a smaller second eigenvalue than a branched structure [55]. Thus, the RNAs are analyzed by means of their graph invariants, which are eigenvalues here.

Labeling the tree graphs with IDs

We label all tree graphs of the same vertex number by increasing λ_2 . Thus, for example, ID 6_1 indicates that the graph has 6 nodes and the smallest λ_2 among all 6-node graphs; ID 6_3 indicates the 6-node graph with the third lowest λ_2 and so on.

Deduction of characteristic information from the Laplacian spectra

To derive essential topological features of an RNA graph so we can compare and visualize, in 2D or 3D, the graphs

with varying number of nodes, we compress the number of descriptors from the Laplacian spectrum, which is composed of n eigenvalues for a graph of n vertices, to two variables α and β : the slope α and the intercept β are calculated by applying the linear least-square regression to the set of planar points $(1, \lambda_2), (2, \lambda_3), \dots, (n-1, \lambda_n)$. The first eigenvalue λ_1 is omitted because its value is always zero. Thus, α measures the average spacing between positive eigenvalues and the intercept β represents the second smallest eigenvalue calibrated by α . This type of reduction mechanism is commonly used in clustering analysis. One example is in the field of drug design, known as quantitative structure–activity relationships [56], where various chemical compounds are described by a few “topological descriptors”.

Here, we observe that α decreases with n , and therefore, we assume that $n\alpha$ forms a quantity independent of n . We thus derive a set of two descriptors, $(n\alpha, \beta)$, and use this quantity as a component to perform clustering of RNA-like and non-RNA-like motifs based on the existing RNA databases. In addition, considering the relationship of the eigenequation for powers $k = 0, 1, 2, \dots$,

$$L^k x_i = \lambda_i^k x_i \quad (i = 1, 2, \dots, n),$$

where x_i is an eigenvector corresponding to λ_i , enhances the accuracy of clustering effectively [15] by allowing us to add more parameters. We define α_k and β_k in the same manner from the powers of the eigenvalues $(1, \lambda_2^k), (2, \lambda_3^k), \dots, (n-1, \lambda_n^k)$. Thus, a point in a $2k$ dimensional space is obtained for each secondary structure. Our previous work [15] showed some advantage of the $k = 2$ space over other values; thus, this value is consistently used here too.

To make each coordinate's contribution equal for the predictions, we normalized these values based on the average of their absolute values. That is, if we let $x_m = (m\text{th coordinate})$, for example, $x_1 = n\alpha_1$, the normalized coordinates x_m^* are

$$x_m^* = (\bar{x}_1 / \bar{x}_m) x_m.$$

Note that, although we chose (\bar{x}_1) for the numerator, this could be the mean of any x_m .

Finally, the metric multidimensional scaling (MDS) is performed to map these four dimensional points to the same number of two dimensional points keeping the Euclidean distances among the original points as much as possible [57].

Program implementation

As mentioned, the 2D tree graphs are generated by the combination of nauty and Traces [38] and NetworkX [53]. The code for converting RNA 2D full topology into a tree graph, which was described in the section [RNA tree graph representation](#), was automated in our previous work [19] and is used here too. The MDS is performed by the implementation of the function `cmdscales` from the multi-variable analysis library package of R [58]. The k -NN and PAM clustering are performed by the C clustering library [20]. All other parts are coded by the first author using Python. The entire calculation process takes less than 2 h on Intel® Core™ i5-4258U.

Acknowledgements

This work is supported by the National Science Foundation (DMS-0201160 and CCF-0727001) and the National Institutes of Health (GM100469 and GM081410). N.B. also would like to thank Professor Irlé at Nagoya University for his support. We also thank Cigdem S. Bayrak for her assistance.

Appendix A. Supplementary data

Supplementary data to this article can be found online at <http://dx.doi.org/10.1016/j.jmb.2015.10.009>.

Received 11 September 2015;

Accepted 6 October 2015

Available online 22 October 2015

Keywords:

RNA secondary structure;
RNA atlas;
RNA motifs;
RNA design;
Prediction of RNA-like motifs

† <http://www.biomath.nyu.edu/rag/home>.

Abbreviations used:

PDB, Protein Data Bank; k -NN, k -nearest neighbor; MDS, multidimensional scaling; PAM, partitioning around medoids.

References

- [1] S.R. Eddy, Non-coding RNA genes and the modern RNA world, *Nat. Rev. Genet.* 2 (2001) 919–929.
- [2] E. Nudler, Flipping riboswitches, *Cell* 126 (2006) 19–22.
- [3] R.R. Breaker, Riboswitches and the RNA world, *Cold Spring Harbor Perspect. Biol.* 4 (2012) a003566.
- [4] S. Gribaldo, C. Brochier-Armanet, The origin and evolution of Archaea: A state of the art, *Philos. Trans. R. Soc. Lond. B Biol. Sci.* 361 (2006) 1007–1022.
- [5] S.K. Burley, et al., Structural genomics: Beyond the Human Genome Project, *Nat. Genet.* 23 (1999) 151–157.
- [6] M.R. Chance, et al., Structural genomics: A pipeline for providing structures for the biologist, *Protein Sci.* 11 (2002) 723–738.
- [7] D. Bray, Molecular networks: The top-down view, *Science* 301 (2003) 1864–1865.
- [8] S. Kalir, U. Alon, Using a quantitative blueprint to reprogram the dynamics of the flagella gene network, *Cell* 117 (2004) 713–720.
- [9] A.L. Barabási, E. Bonabeau, Scale-free networks, *Sci. Am.* 288 (2003) 60–69.
- [10] S.H. Yook, H. Jeong, A.L. Barabási, Modeling the Internet's large-scale topology, *Proc. Natl. Acad. Sci. U. S. A.* 99 (2002) 13382–13386.
- [11] P. Santi, et al., Quantifying the benefits of vehicle pooling with shareability networks, *Proc. Natl. Acad. Sci. U. S. A.* 111 (2014) 13290–13294.
- [12] S.Y. Le, R. Nussinov, J.V. Maizel, Tree graphs of RNA secondary structures and their comparisons, *Comput. Biomed. Res.* 22 (1989) 461–473.
- [13] G. Benedetti, S. Morosetti, A graph-topological approach to recognition of pattern and similarity in RNA secondary structures, *Biophys. Chem.* 59 (1996) 179–184.
- [14] W. Fontana, D.A.M. Konings, P.F. Stadler, P. Schuster, Statistics of RNA secondary structures, *Biopolymers* 33 (1993) 1389–1404.
- [15] N. Kim, N. Shiffeldrim, H.H. Gan, T. Schlick, Candidates for novel RNA topologies, *J. Mol. Biol.* 341 (2004) 1129–1144.
- [16] N. Kim, N. Fuhr, T. Schlick, Graph Applications to RNA Structure and Function Chapter 3 in: R. Russell (Ed.), *Biophysics of RNA Folding, Biophysics for the Life Sciences*, 3, Springer Verlag 2013, pp. 23–51.
- [17] N. Kim, L. Petingi, T. Schlick, Network theory tools for RNA modeling, *WSEAS Trans. Acoust. Math* 12 (2013) 941–955.
- [18] N. Kim, M. Zahran, T. Schlick, Computational prediction of riboswitch tertiary structures including pseudoknots by RAGTOP: A hierarchical graph sampling approach, *Methods Enzymol.* 553 (2015) 115–135.
- [19] J.A. Izzo, N. Kim, S. Elmetwaly, T. Schlick, RAG: An update to the RNA-As-Graphs resource, *BMC Bioinf.* 12 (2011) 219.
- [20] M. de Hoon, S. Imoto, S. Miyano, The C Clustering Library, Institute of Medical Science Human Genome Center, University of Tokyo, Tokyo, Japan, 2005.
- [21] T. Haynes, D. Knisley, E. Seier, Y. Zou, A quantitative analysis of secondary RNA structure using domination based parameters on trees, *BMC Bioinf.* 7 (2006) 108.
- [22] D.R. Koessler, D.J. Knisley, J. Knisley, T. Haynes, A predictive model for secondary RNA structure using graph theory and a neural network, *BMC Bioinf.* 11 (2010) S21.
- [23] M. Hamada, K. Tsuda, T. Kudo, T. Kin, K. Asai, Mining frequent stem patterns from unaligned RNA sequences, *Bioinformatics* 22 (2006) 2480–2487.
- [24] U. Laserson, H.H. Gan, T. Schlick, Predicting candidate genomic sequences that correspond to synthetic functional RNA motifs, *Nucleic Acids Res.* 33 (2005) 6057–6069.
- [25] N. Kim, J.S. Shin, S. Elmetwaly, H.H. Gan, T. Schlick, RAGPools: RNA-As-Graph-Pools—A Web server for assisting the design of structured RNA pools for *in vitro* selection, *Bioinformatics* 23 (2007) 2959–2960.
- [26] N. Kim, J.A. Izzo, S. Elmetwaly, H.H. Gan, T. Schlick, Computational generation and screening of RNA motifs in large nucleotide sequence pools, *Nucleic Acids Res.* 38 (2010) e139.
- [27] N. Kim, H.H. Gan, T. Schlick, A computational proposal for designing structured RNA pools for *in vitro* selection of RNAs, *RNA* 13 (2007) 478–492.
- [28] A. Gopal, Z.H. Zhou, C.M. Knobler, Visualizing large RNA molecules in solution, *RNA* 18 (2012) 284–299.
- [29] Y. Bakhtin, C.E. Heitsch, Large deviations for random trees and the branching of RNA secondary structures, *Bull. Math. Biol.* 71 (2009) 84–106.
- [30] G. Quarta, N. Kim, J.A. Izzo, T. Schlick, Analysis of riboswitch structure and function by an energy landscape framework, *J. Mol. Biol.* 393 (2009) 993–1003.
- [31] G. Quarta, K. Sin, T. Schlick, Dynamic energy landscapes of riboswitches help interpret conformational rearrangements and function, *PLoS Comput. Biol.* 8 (2012) e1002368.

- [32] N. Kim, Z. Zheng, S. Elmetwaly, T. Schlick, RNA graph partitioning for the discovery of RNA modularity: A novel application of graph partition algorithm to biology, *PLoS One* 9 (2014) e106074.
- [33] N. Kim, C. Laing, S. Elmetwaly, S. Jung, J. Curuksu, T. Schlick, Graph-based sampling for approximating global helical topologies of RNA, *Proc. Natl. Acad. Sci. U. S. A.* 111 (2013) 4079–4084.
- [34] C. Laing, S. Jung, N. Kim, S. Elmetwaly, M. Zahran, T. Schlick, Predicting helical topologies in RNA junctions as tree graphs, *PLoS One* 8 (2013) e71947.
- [35] P.P. Gardner, et al., Rfam: Updates to the RNA families database, *Nucleic Acids Res.* 37 (2009) D136–D140.
- [36] S.G. Jones, A. Bateman, M. Marshall, A. Khanna, S.R. Eddy, Rfam: An RNA family database, *Nucleic Acids Res.* 31 (2003) 439–441.
- [37] M. Andronescu, V. Bereg, H.H. Hoos, A. Condon, RNA STRAND: The RNA secondary structure and statistical analysis database, *BMC Bioinf.* 9 (2008) 340.
- [38] B.D. McKay, A. Piperno, Practical graph isomorphism, II, *J. Symb. Comput.* 60 (2013) 94–112.
- [39] L. Kaufman, P.J. Rousseeuw, *Finding Groups in Data: An Introduction to Cluster Analysis*, Wiley Interscience, Hoboken, 1990.
- [40] B.D. Ripley, *Pattern Recognition and Neural Networks*, Cambridge University Press, Cambridge, 1996.
- [41] W.N. Venables, B.D. Ripley, *Modern Applied Statistics with S*, fourth ed. Springer, New York, 2002.
- [42] L. Torgo, *Data Mining with R: Learning with Case Studies*, Chapman & Hall/CRC Boca Raton, 2011.
- [43] A.T. Larios, K.K. Swinger, A.S. Krasilnikov, T. Pan, A. Mondragon, Crystal structure of the RNA component of bacterial ribonuclease P, *Nature* 437 (2005) 584–587.
- [44] J.E. Johnson, F.E. Reyes, J.T. Polaski, R.T. Batey, B12 cofactors directly stabilize an mRNA regulatory switch, *Nature* 492 (2012) 133–137.
- [45] F. Weis, et al., tmRNA-SmpB: A journey to the centre of the bacterial ribosome, *EMBO J.* 29 (2010) 3810–3818.
- [46] E.P. Nawrocki, S.W. Burge, A. Bateman, J. Daub, R.Y. Eberhardt, S.R. Eddy, E.W. Floden, P.P. Gardner, T.A. Jones, J. Tate, R.D. Finn, Rfam 12.0: Updates to the RNA families database, *Nucleic Acids Res.* 43 (2015) D130–D137, <http://dx.doi.org/10.1093/nar/gku1063>.
- [47] F.H.D. van Batenburg, A.P. Gulyaev, C.W.A. Pleij, J. Ng, J. Oliehoek, PseudoBase: A database with RNA pseudoknots, *Nucleic Acids Res.* 28 (2000) 201–204.
- [48] M. Andronescu, V. Bereg, H.H. Hoos, A. Condon, RNA STRAND: The RNA secondary structure and statistical analysis database, *BMC Bioinf.* 13 (2008) 340.
- [49] H.M. Berman, et al., The Protein Data Bank, *Nucleic Acids Res.* 28 (2000) 235–242.
- [50] H.M. Berman, et al., The Nucleic Acid Database: A comprehensive relational database of three-dimensional structures of nucleic acids, *Biophys. J.* 63 (1992) 751–759.
- [51] B.C. Narayanan, et al., The Nucleic Acid Database: New features and capabilities, *Nucleic Acids Res.* 42 (2013) D114–D122.
- [52] F. Harary, *Graph Theory*, Perseus Books, Reading, 1999.
- [53] A.A. Hagberg, D.A. Schult, P.J. Swart, *Exploring Network Structure, Dynamics, and Function Using NetworkX*, Proceedings of the 7th Python in Science Conference, 2008 (SciPy2008).
- [54] A.E. Brouwer, W.H. Haemers, *Spectra of Graphs*, Springer, New York, 2012.
- [55] F. Chung, *Spectral Graph Theory*, Published for the Conference Board of the Mathematical Sciences by the American Mathematical Society, Providence, 1997.
- [56] T. Schlick, *Molecular Modeling and Simulation an Interdisciplinary Guide*, Springer, New York, 2002.
- [57] I. Borg, P.J.F. Groenen, *Modern Multidimensional Scaling Theory and Applications*, second ed. Springer, New York, 2005.
- [58] R Development Core Team, *R: A Language and Environment for Statistical Computing*, R Foundation for Statistical Computing, Vienna, Austria, 2011.