

BPPart and BPMax: RNA-RNA Interaction Partition Function and Structure Prediction for the Base Pair Counting Model

Ali Ebrahimpour-Boroojeny, Sanjay Rajopadhye, and Hamidreza Chitsaz *

Department of Computer Science, Colorado State University

Abstract

A few elite classes of RNA-RNA interaction (RRI), with complex roles in cellular functions such as miRNA-target and lncRNAs in human health, have already been studied. Accordingly, RRI bioinformatics tools tailored for those elite classes have been proposed in the last decade. Interestingly, there are somewhat unnoticed mRNA-mRNA interactions in the literature with potentially drastic biological roles. Hence, there is a need for high-throughput *generic* RRI bioinformatics tools.

We revisit our RRI partition function algorithm, **piRNA**, which happens to be the most comprehensive and computationally-intensive thermodynamic model for RRI. We propose simpler models that are shown to retain the vast majority of the thermodynamic information that **piRNA** captures.

We simplify the energy model and instead consider only weighted base pair counting to obtain **BPPart** for Base-pair Partition function and **BPMax** for Base-pair Maximization which are $225\times$ and $1350\times$ faster than **piRNA**, with a correlation of 0.855 and 0.836 with **piRNA** at $37^\circ C$ on 50,500 experimentally characterized RRIs. This correlation increases to 0.920 and 0.904, respectively, at $-180^\circ C$.

Finally, we apply our algorithm **BPPart** to discover two disease-related RNAs, SNORD3D and TRAF3, and hypothesize their potential roles in Parkinson’s disease and Cerebral Autosomal Dominant Arteriopathy with Subcortical Infarcts and Leukoencephalopathy (CADASIL).

1 Introduction

Since mid 1990s with the advent of RNA interference discovery, RNA-RNA interaction (RRI) has moved to the spotlight in modern, post-genome biology. RRI is ubiquitous and has increasingly complex roles in cellular functions. In human health studies, miRNA-target and lncRNAs are among an elite class of RRIs that have been extensively studied and shown to play significant roles in various diseases including cancer. Bacterial ncRNA-target and RNA interference are other classes of RRIs that have received significant attention. However, new evidence suggests that other classes of RRI, such as mRNA-mRNA interactions, are biologically important.

The RISE database [1] reports a number of biologically significant instances of mRNA-mRNA interactions. These representative mRNA-mRNA interactions suggest that general RRIs, including mRNA-mRNA interactions, play major roles in human biology. Hence, there is a need for high-throughput *generic* RNA-RNA interaction bioinformatic tools for all types of RNAs.

As an example of this necessity for all types of RNAs, we found 3 cliques of size 4 of interacting protein-coding RNAs in ribosome which conform to what we generally expect from the structure of the ribosome. These cliques are highly entangled together to form an interaction graph as Figure 1. RPS3 which seems to be one of the genes with the highest number of connections interacts with at least 14 other genes in ribosome pathway. Another interesting clique of size 4 that we could find consists of 4 genes in the pathway of regulation of actin cytoskeleton, ACTB, ACTG1, PFN1, and TMSB4X. These genes are involved in vital tasks of proliferation, migration, mobility, and differentiation of the cell. Being able to capture all the

*To whom correspondence should be addressed. Email: chitsaz@chitsazlab.org

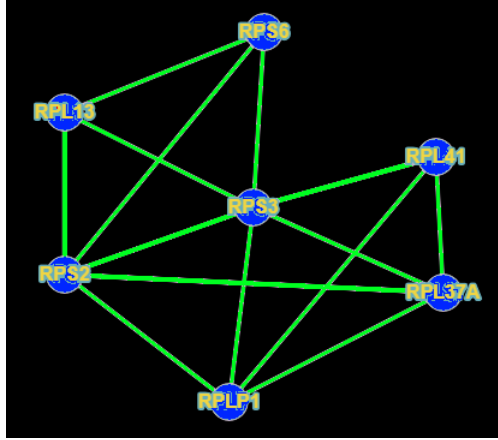


Figure 1: A substructure of the genes in the ribosome pathway. Each node represents a gene and each edge represents an experimentally observed interaction between the corresponding genes.

interactions that RNAs might have will help us better understand the post-transcriptional regulation of the genes.

In this paper, we revisit our RNA-RNA interaction partition function algorithm, **piRNA**, which happens to be the most comprehensive, albeit the most computationally intensive, thermodynamic model for RNA-RNA interaction [2]. **piRNA** is a dynamic programming algorithm that computes the partition function, base-pairing probabilities, and structure for the comprehensive Turner energy model in $O(n^4m^2 + n^2m^4)$ time and $O(n^4 + m^4)$ space. Due to intricacies of the energy model, including various (kissing) loops such as hairpin loop, bulge/internal loop, and multibranch loop, **piRNA** involves 96 different dynamic programming tables and needs multiple table look-ups for computing their values. An implementation of **piRNA** is currently available at <http://chitsazlab.org/software/piRNA>.

In this paper, we introduce a strategic retreat from the slower comprehensive models such as **piRNA** by simplifying the energy model and instead considering only simple weighted base pair counting. We develop the **BPPart** algorithm for Base-pair Partition function and **BPMMax** for Base-pair Maximization, both of which are faster by a significant, albeit constant factor. By the explosion of experimental data which makes us able to use machine learning methods, such as deep learning, for detection of RNA subsequences that interact, this retreat is necessary if one is willing to build physics-guided models by using the features that are derived by an energy model. **BPPart** involves nine 4-dimensional dynamic programming tables, and **BPMMax** involves only one 4-dimensional table. Both **BPPart** and **BPMMax** compared with **piRNA** are simpler dynamic programming algorithms which are more than $225\times$ and $1300\times$ faster, respectively, on the 50,500 RRI samples we used for our experiments. The reason for this noticeable speed-up is reducing the number of tables and the number of table look-ups for computing the new values and also the fact that the 96 large tables of **piRNA** renders **piRNA** memory- rather than compute-bound in practice. Moreover, the significantly reduced memory footprint of **BPPart** and **BPMMax** makes them feasible targets for optimization on different hardware platforms like GPU based accelerators, an avenue we plan to explore in the future.

The key question concerns the accuracy we lose by simplifying the scoring model from the comprehensive Turner model to simply weighted base-pair counting. We answer this by computing both the Pearson and Spearman’s rank correlations at different temperatures between the results of **BPPart**, **BPMMax**, and **piRNA** on 50,500 experimentally characterized RRIs in the RISE database [1]. We find that the Pearson correlations between **BPPart** and **piRNA** is 0.920 and **BPMMax** and **piRNA** is 0.904 at -180°C after optimizing the weights for base pairs. The effect of entropy, for which the simple base pair counting model does not account, increases with temperature. Completely in conformance with this theoretical expectation, we find that the Pearson correlations between **BPPart** and **piRNA** is 0.855 and **BPMMax** and **piRNA** is 0.836 at 37°C . We conclude that **BPPart** and **BPMMax** capture a significant portion of the thermodynamic information. They can possibly be complemented with machine learning techniques in the future for more accurate predictions.

1.1 Related work

During the last few decades, several computational methods emerged to study the secondary structure of single and interacting nucleic acid strands. Most use a thermodynamic model such as the well-known Nearest Neighbor Thermodynamic model [3, 4, 5, 2, 6, 7, 8, 9, 10, 11]. Some previous attempts to analyze the thermodynamics of multiple interacting strands concatenate input sequences in some order and consider them as a single strand [12, 13, 14]. Alternatively, several methods avoid internal base-pairing in either strand and compute the minimum free energy secondary structure for their hybridization under this constraint [15, 16, 17]. The most comprehensive solution is computing the joint structure between two interacting strands under energy models with a growing complexity [18, 19, 20, 21, 22, 2, 23].

Other methods predict the secondary structure of individual RNA independently, and predict the (most likely) hybridization between the unpaired regions of the two interacting molecules as a multistep process: 1) unfolding of the two molecules to expose bases needed for hybridization, 2) the hybridization at the binding site, and 3) restructuring of the complex to a new minimum free energy conformation [24, 25, 26, 27]. The success of such methods, including our biRNA algorithm [27], suggests that the thermodynamic information vested in subsequences and pairs of subsequences of the input RNAs can provide valuable information for predicting features of the entire interaction.

In addition to general RNA-RNA interaction tools, many tools have been developed to predict the secondary structure of interacting RNAs for a specific type of interest which has been shown to be more effective in some cases due to the utilization of certain properties belonging to that type. As mentioned earlier, miRNA-target prediction is one such class of high interest for which such specialized tools have been created to incorporate various properties specific to miRNAs; some of these tools use the seed region of a miRNA which is highly conserved [28, 29, 30, 31], some consider the free energy to compute accessibility to the binding site in 3' UTR [32, 20, 29], some utilize the conservation level which is derived using the phylogenetic distance [33, 34, 35, 36, 28, 29], and some others consider other target sites as well, such as the 5' UTR, Open Reading Frames (ORF), and the coding sequence (CDS) for mRNAs [37, 38, 39, 40].

There are also several other tools developed for other specific types of RNA; IntaRNA [41, 42] is one such tool that although is used for RNA-RNA interaction in general, it is primarily designed for predicting target sites of non-coding RNAs (ncRNAs) on mRNAs. There are many other examples, such as PLEXY [43] which is a tool designed for C/D snoRNAs, RNAsnoop [44] that is designed for H/ACA snoRNAs, TargetRNA [45] which is a tool aimed at predicting interaction of bacterial sRNAs [46].

2 Methods

Here we describe how our algorithm, **BPPart**, utilizes a dynamic programming approach to compute the partition function for RNA-RNA interaction when entropy is ignored and only a weighted score for pairing different nucleotides is considered. This algorithm is guaranteed to be mutually exclusive on the set of structures, i.e., it counts each structure exactly once. For **BPMax** which maximizes the (weighted scores) of base-pairs, such mutual exclusion is not necessary because the max operator is idempotent (counting the same structure multiple times does not affect the value of the objective function) and we give a $10\times$ simpler recursion. Our codes are freely available under open source license.¹

Preliminaries

In this paper, we mostly follow the notations and definitions used to develop our **piRNA** algorithm [2]. We denote the two nucleic acid strands by **R** and **S**. Strand **R** is indexed from 1 to L_R , and **S** is indexed from 1 to L_S both in 5' to 3' direction. Note that the two strands interact in opposite directions, e.g. **R** in $5' \rightarrow 3'$ with **S** in $3' \leftarrow 5'$ direction; however, we consider the reverse of **S** in our figures for clearer illustration of the configurations. Each nucleotide is paired with at most one nucleotide in the same or the other strand. The subsequence from the i^{th} nucleotide to the j^{th} nucleotide, inclusive, in either strand is denoted by $[i, j]$.

An intramolecular base pair between the nucleotides i and j (by convention, $i < j$) in a strand is called an *arc* and denoted by a bullet $i \bullet j$. We represent the score of such arc by $\text{score}(i, j)$. Essentially, $\text{score}(i, j)$

¹See <https://github.com/Ali-E/bipart>

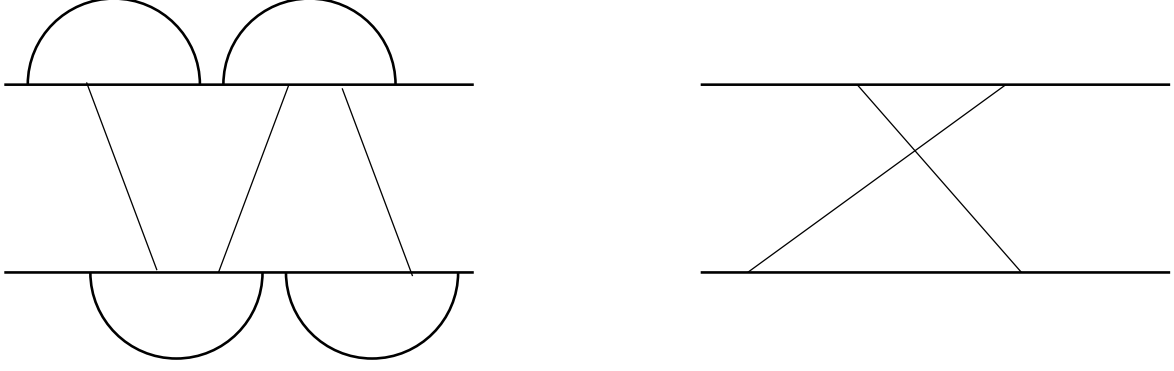


Figure 2: An illustration of a zigzag (left) and a crossing bond (right), which are excluded in our algorithm.

is c_1 if $i \bullet j$ is GU or UG, is c_2 if $i \bullet j$ is AU or UA, and is c_3 if $i \bullet j$ is CG or GC. An intermolecular base pair between the nucleotides k_1 and k_2 , where $k_1 \in R, k_2 \in S$, is called a *bond*, denoted by a circle $k_1 \circ k_2$. We represent the score of such a bond by $\text{iscore}(k_1, k_2)$. Essentially, $\text{iscore}(k_1, k_2)$ is c'_1 if $k_1 \circ k_2$ is GU or UG, is c'_2 if $k_1 \circ k_2$ is AU or UA, and is c'_3 if $k_1 \circ k_2$ is CG or GC.

An arc $i \bullet j$ in R *covers* a bond $k_1 \circ k_2$ if $i_1 < k_1 < j_1$. We call $i \bullet j$ an *interaction arc* in R if there is a bond $k_1 \circ k_2$ covered by $i \bullet j$. The *scope* of an interaction arc is the interval $[i + 1, j - 1]$. We call a base on either strand an *event* if it is either the end point of a bond or that of an interaction arc. In our explanation we may use arc and bond as verbs. Two bonds $i_1 \circ i_2$ and $j_1 \circ j_2$ are called *crossing bonds* (right case of Figure 2) if $i_1 < j_1$ and $i_2 > j_2$, or vice versa. An interaction arc $i_1 \bullet j_1$ in a strand *subsumes* a subsequence $[i_2, j_2]$ in the other strand if none of the bases in $[i_2, j_2]$ has a bond with a base outside this arc. Mathematically, for all bonds $k_1 \circ k_2$ where $i_2 < k_2 < j_2$, k_1 lies within the scope of $i_1 \bullet j_1$. Two interaction arcs are *equivalent* if they subsume one another. Two interaction arcs $i_1 \bullet j_1$ and $i_2 \bullet j_2$ are part of a *zigzag*, if neither $i_1 \bullet j_1$ subsumes $[i_2, j_2]$ nor $i_2 \bullet j_2$ subsumes $[i_1, j_1]$ (left case of Figure 2).

In this work, we assume there are no pseudoknots in individual secondary structures of R and S , and also there are no crossing bonds and no zigzags between R and S . These constraints allow a polynomial algorithm—the general case of considering all possible structures is NP-hard [19]. We denote the ensemble of unspseudoknotted structures of R and S by $\mathcal{S}(R)$ and $\mathcal{S}(S)$ respectively. The ensemble of unspseudoknotted, crossing-free, and zigzag-free joint interaction structures is denoted by $\mathcal{S}^I(R, S)$.

For a given structure s in either $\mathcal{S}(R)$ or $\mathcal{S}(S)$, let $\text{AU}(s)$ denote the number of A-U base pairs in s . Similarly, $\text{CG}(s)$ and $\text{GU}(s)$ denote the number of C-G and G-U base pairs in s , respectively. We define *bpcount* as a weighted sum, for some constants, c_1, \dots, c_3

$$\text{bpcount}(s) = c_1 \text{GU}(s) + c_2 \text{AU}(s) + c_3 \text{CG}(s). \quad (1)$$

For a given joint interaction structure $s \in \mathcal{S}^I(R, S)$, let $\text{AU}(s)$, $\text{CG}(s)$, and $\text{GU}(s)$ denote the respective number of intramolecular base pairs in s , and let $\text{AU}^I(s)$, $\text{CG}^I(s)$, and $\text{GU}^I(s)$ denote the number of corresponding intermolecular base pairs in s . We define for some constants, c'_1, \dots, c'_3 , for any joint interaction structure s ,

$$\text{bpcount}^I(s) = c'_1 \text{GU}^I(s) + c'_2 \text{AU}^I(s) + c'_3 \text{CG}^I(s), \quad (2)$$

and

$$\text{bpcount}(s) = c_1 \text{GU}(s) + c_2 \text{AU}(s) + c_3 \text{CG}(s) + \text{bpcount}^I(s). \quad (3)$$

Rivas-Eddy Diagrams

For the sake of completeness, we describe the ‘‘Rivas-Eddy diagram’’ notation that we adopt in this paper. The main elements are:

1. A solid horizontal straight line represents a sequence; we have two sequences drawn as two parallel horizontal lines.
2. A solid curved line between two points in the same sequence is an arc; all arcs are either above the upper sequence, or below the lower one.
3. A dotted curved line with a cross in the middle, between two points in the same sequence means that those two points *do not* form an arc.
4. A dashed curved line between two points in the same sequence denotes either 2 or 3.
5. A solid line between two points in different sequences is a bond.
6. Similarly, a dotted line with a cross in the middle, between two points in different sequences means that those two points *do not* form a bond.
7. A dashed line between two points in different sequences denotes either 5 or 6.
8. A region is the space under an arc, or between bonds. When there are no additional choices of bonds/arcs in a given region, we fill it with a color (cyan); no arc or bond crosses a filled region.
9. A point in a sequence may be labeled with an index, and in general, the set of such indices are free variables used in the recursions; the index of unlabeled points before (after) labeled points is assumed to be the predecessor (successor) of the label.
10. A diagram may be labeled with the name(s) of the constituent substructures (which are eventually implemented as dynamic programming tables/variables).
11. A vanishing arc (i.e., one that starts at some index, and does not explicitly specify an end point) represents a structure whose start point is as specified, and the end point is to be determined.

The Rivas-Eddy diagram to compute a certain function is written like a formal (context free) grammar. The left hand side is labeled with the name of a table (structure), and the right hand side has a number of alternate substructures separated by vertical bars. Often, some of the boundary cases (e.g., singleton or empty subsequences) are omitted for brevity.

Problem Definition

In this paper, we solve two problems:

1. **Base Pair Counting Partition Function:** we give a dynamic programming algorithm **BPPart** to compute the partition function

$$Q(\mathbf{R}, \mathbf{S}) = \sum_{s \in S^I(\mathbf{R}, \mathbf{S})} e^{\text{bpcount}(s)}, \quad (4)$$

2. **Base Pair Maximization:** we give a dynamic programming algorithm **BPMax** to find the structure that has the maximum weighted base pair count, i.e.

$$\text{BPMax}(\mathbf{R}, \mathbf{S}) = \max_{s \in S^I(\mathbf{R}, \mathbf{S})} \text{bpcount}(s). \quad (5)$$

This problem was previously studied by Pervouchine [18] in an algorithm called IRIS. However, there is no publicly available correct implementation of IRIS. Moreover, we also define an additional interaction score to capture the structure with the highest intermolecular score, amongst those that maximize the total score. Mathematically,

$$\text{IS}(\mathbf{R}, \mathbf{S}) = \max_{\{s \mid \text{bpcount}(s) = \text{BPMax}(\mathbf{R}, \mathbf{S})\}} \text{bpcount}^I(s). \quad (6)$$

We compute $\text{IS}(\mathbf{R}, \mathbf{S})$ by backtracing all possible total-score-optimal structures, and selecting the one that has the maximum intermolecular score.

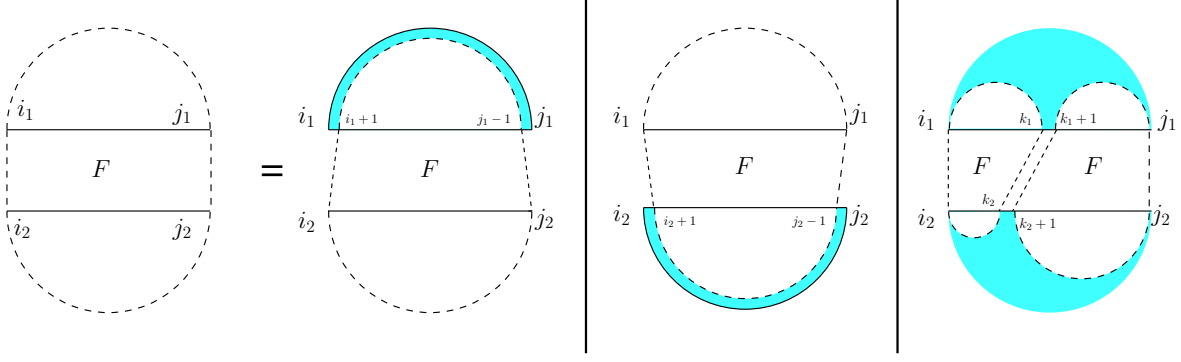


Figure 3: The four cases defining table F . Note that in the **BPM** algorithm, the cases do not have to be mutually exclusive since we are working with the max operator, which is idempotent.

BPM

We first explain the **BPM** algorithm. It is simpler than **BPPart**, and allows us to describe the notation and conventions. When explaining some of the equations, helper functions, called H, L, M, N , are used to ease the reading of the paper. To differentiate these helper functions, superscripts are used.

For a single strand of nucleotides, we define $S_{i,j}$ as the maximum weighted sum of base pair scores on all possible foldings of subsequence $[i, j]$. We need to make such a table, for each of the **R** and **S** strands, and we use superscripts (1) and (2), respectively, to distinguish between them. We also define F_{i_1, j_1, i_2, j_2} as the maximum weighted sum of base pair scores (considering both intra- and inter-pairings) of subsequences $[i_1, j_1]$ from **R** and $[i_2, j_2]$ from **S**.

The computation of $S_{i,j}$ is based on the well known single RNA folding algorithm [47]. For short sequences (i.e., those whose length is strictly less than 5) the score is 0, otherwise, we use the recursion in the second case of Equation (7) shown below. It considers the case where we have an arc $i \bullet j$ and recurs on $[i+1, j-1]$, and also other cases in which the i^{th} and j^{th} bases are not paired and the $[i, j]$ is split into two smaller subsequences:

$$S_{i,j} = \begin{cases} 0 & j - i < 4 \\ \max \left(S_{i+1, j-1} + \text{score}(i, j), \max_{k=i}^{j-1} S_{i,k} + S_{k+1, j} \right) & \text{otherwise.} \end{cases} \quad (7)$$

We define the recurrences for F_{i_1, j_1, i_2, j_2} similarly. When either sequence is empty, the value is simply the S of the other sequence, and for two singleton sequences, it is the score of the single bond possible. Otherwise we have three cases: (i) i_1 arcs j_1 ($i_1 \bullet j_1$) in which case the residual structure is given by a recursion on $F_{i_1+1, j_1-1, i_2, j_2}$, (ii) the symmetric case of $i_2 \bullet j_2$ and $F_{i_1, j_1, i_2+1, j_2-1}$, or (iii) none of these arcs, and two recursive cases of F_{i_1, k_1, i_2, k_2} and $F_{k_1+1, j_1, k_2+1, j_2}$. They are illustrated in the Rivas-Eddy diagram of Figure 3, which lead to

$$F_{i_1, j_1, i_2, j_2} = \begin{cases} -\infty & j_1 < i_1 \text{ and } j_2 < i_2 \\ S_{i_1, j_1}^{(1)} & i_1 \leq j_1 \text{ and } j_2 < i_2 \\ S_{i_2, j_2}^{(2)} & j_1 < i_1 \text{ and } i_2 \leq j_2 \\ \text{iscore}(i_1, i_2) & i_1 = j_1 \text{ and } i_2 = j_2 \\ \max [F_{i_1+1, j_1-1, i_2, j_2} + \text{score}(i_1, j_1), \\ F_{i_1, j_1, i_2+1, j_2-1} + \text{score}(i_2, j_2), \\ H_{i_1, j_1, i_2, j_2}] & \text{otherwise,} \end{cases} \quad (8)$$

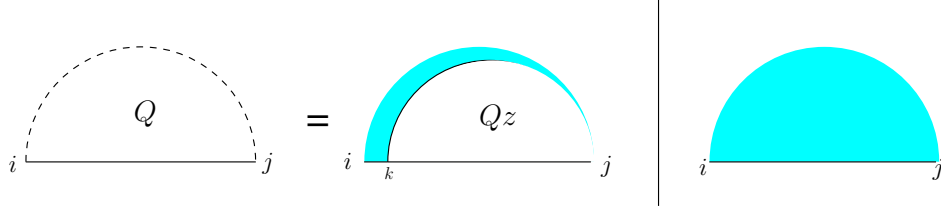


Figure 4: For computing Q , notice that either there is no pairing or there is at least one arc which starts at some index k and results in a case of Qz .

$$H_{i_1, j_1, i_2, j_2} = \max_{k_1=i_1-1}^{j_1} \max_{k_2=i_2-1}^{j_2} (F_{i_1, k_1, i_2, k_2} + F_{k_1+1, j_1, k_2+1, j_2}). \quad (9)$$

Note that H is equivalent to

$$H_{i_1, j_1, i_2, j_2} = \max \left(\begin{array}{l} S^{(1)}(i_1, j_1) + S^{(2)}(i_2, j_2), \\ \max_{k_1=i_1}^{j_1-1} \max_{k_2=i_2}^{j_2-1} F_{i_1, k_1, i_2, k_2} + F_{k_1+1, j_1, k_2+1, j_2}, \\ \max_{k_2=i_2}^{j_2-1} S^{(2)}(i_2, k_2) + F_{i_1, j_1, k_2+1, j_2}, \\ \max_{k_2=i_2}^{j_2-1} F_{i_1, j_1, i_2, k_2} + S^{(2)}(k_2+1, j_2), \\ \max_{k_1=i_1}^{j_1-1} S^{(1)}(i_1, k_1) + F_{k_1+1, j_1, i_2, j_2}, \\ \max_{k_1=i_1}^{j_1-1} F_{i_1, k_1, i_2, j_2} + S^{(1)}(k_1+1, j_1) \end{array} \right). \quad (10)$$

In Equation (8), we compute S tables separately for each strand, according to Equation (7) with the corresponding sequence as the input, and we distinguish them by superscripts ⁽¹⁾ and ⁽²⁾ above. We use the same superscript convention throughout this paper.

BPPart Algorithm

It is well known that the partition function can be computed by developing similar recursions, with two simple modifications. The first is that algebraically, we operate with the field of reals rather than the max-plus semi-ring. Here, the additive identity is 0, rather than INT_MIN and the multiplicative identity is 1, rather than 0. The second, as already mentioned earlier, is that because addition is not idempotent, we must carefully ensure that we enumerate substructures in a mutually exclusive manner.

First, we start with the recursions for computing the partition function on a single strand which is going to occur in many cases of the double-stranded version. Let $Q_{i,j}$ represent the partition function of the subsequence $[i, j]$. As shown in Figure 4, there are two mutually exclusive cases: either (the right case) there is no arc, or (the left case) there is a unique leftmost arc (the cyan fill ensures this) which starts at k , and a substructure on $[k, j]$ with an arc starting at k , for which we introduce a new table Qz .

To define $Qz_{i,j}$, let $i \bullet k$ (read as “let i arc k ”) for some index k . This results in two Q substructures, one on $[i+1, k-1]$, and the other on $[k+1, j]$. Therefore, the value of $Qz_{i,j}$ can be computed using Equation (12) which accounts for the assumption that no pairing is allowed between two bases that are less than 3 bases apart:

$$Q_{i,j} = \begin{cases} 1 & j \leq i \\ 1 + \sum_{k=i}^{j-4} Qz_{k,j} & \text{otherwise,} \end{cases} \quad (11)$$

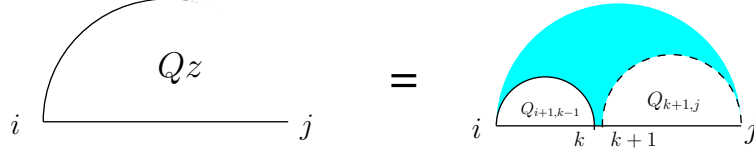


Figure 5: Computing Qz can be achieved by considering the base k that is paired with i and the two Q substructures it forms, one between i and k and one after k .

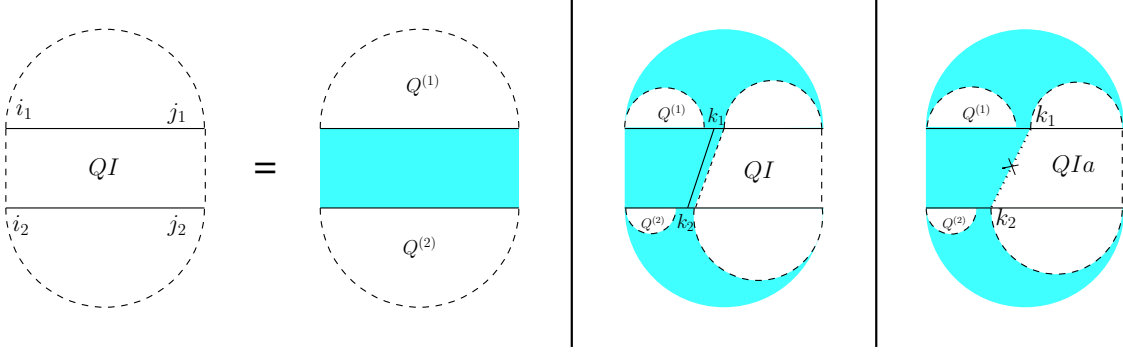


Figure 6: Each case of a QI structure (left side of the equation) can lead to three cases: either no bonds exist (leftmost case), or at least one bond exists. If the first event on both of the sequences is a bond (middle case) the subsequences to the left of the bond involve only Q and the subsequences to the right recurs on QI . Otherwise (rightmost case) we will have QIa (see Figure 7).

$$Qz_{i,j} = \begin{cases} 0 & j - i < 4 \\ \sum_{k=i+4}^j Q_{i+1,k-1} \times e^{\text{score}(i,k)} \times Q_{k+1,j} & \text{otherwise.} \end{cases} \quad (12)$$

For the partition function of a pair of RNA sequences, we consider a 4-dimensional table QI in which QI_{i_1,j_1,i_2,j_2} is the value of base pair counting partition function for the subsequences $[i_1, j_1]$ on \mathbf{R} and $[i_2, j_2]$ on \mathbf{S} . As Figure 6 shows, we can split the set of all possible structures of QI into three mutually exclusive subsets. The leftmost case shows the structures in which there exist no bonds (the first term of Equation (13)). The other two cases occur when there is at least one bond, and hence, unique leftmost events on both \mathbf{R} and \mathbf{S} , at positions k_1 and k_2 , respectively. In the second (middle) case, these leftmost events are end points of a bond, $k_1 \circ k_2$; hence, this case can be broken into: a bond-free section on the left of the bond itself, and a general case of QI on the right of the bond. The third case occurs when k_1 and k_2 are not end points of a bond. We call this structure QIa , and explain it next.

$$\begin{aligned} QI_{i_1,j_1,i_2,j_2} = & Q_{i_1,j_1}^{(1)} \times Q_{i_2,j_2}^{(2)} + \\ & \sum_{k_1=i_1}^{j_1} \sum_{k_2=i_2}^{j_2} L_{i_1,j_1,k_1,i_2,j_2,k_2} + \\ & \sum_{k_1=i_1}^{j_1} \sum_{k_2=i_2}^{j_2} \left(Q_{i_1,k_1-1}^{(1)} \times Q_{i_2,k_2-1}^{(2)} \times QIa_{k_1,j_1,k_2,j_2} \right), \end{aligned} \quad (13)$$

$$L_{i_1,j_1,k_1,i_2,j_2,k_2} = Q_{i_1,k_1-1}^{(1)} \times Q_{i_2,k_2-1}^{(2)} \times e^{\text{score}(k_1,k_2)} \times QI_{k_1+1,j_1,k_2+1,j_2}. \quad (14)$$

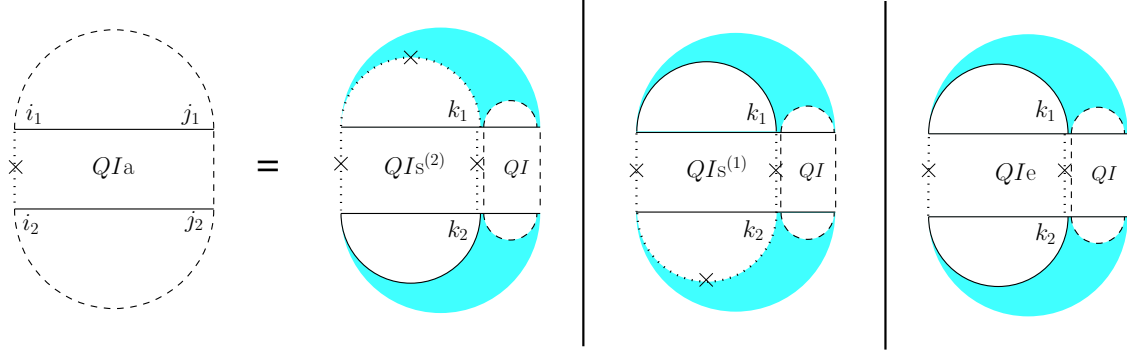


Figure 7: There are three cases for computing the QIa structure; either the leftmost base of only one of the strands is an end point of an arc or both end points are.

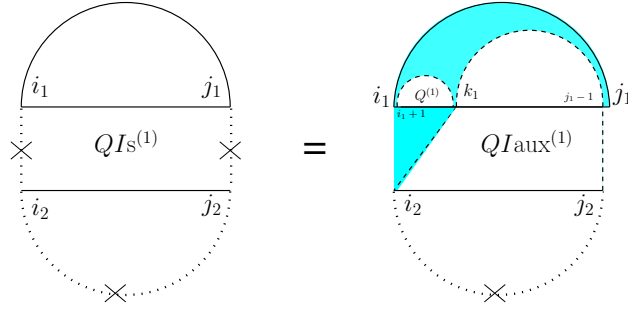


Figure 8: $QIs^{(1)}$ has one arc that can be extracted and the structure derived will have the property that the two end bases of the bottom strand cannot be paired (the new structure inherits this property from $QIs^{(1)}$). On the top strand, we consider the leftmost event. This new structure is $QIaux^{(1)}$.

For computing QIa_{i_1, j_1, i_2, j_2} , (see Figure 7) we have to consider the property of this structure that the leftmost bases on both \mathbf{R} and \mathbf{S} have to be events, but they cannot both be the end points of a bond. Therefore, either one or both of them have to be end points of an interaction arc. There are two possibilities.

First, if both i_1 and i_2 are end points of some interaction arcs, $i_1 \bullet k_1$ and $i_2 \bullet k_2$, these arcs must be equivalent (or else, we have a zigzag). As shown in the rightmost diagram in Figure 7, QIa then splits into two exclusive substructures, namely one where the first and last bases on each strand are paired, and the two arcs are equivalent (we call it QIe_{i_1, k_1, i_2, k_2} and derive its recursion later), followed by $QI_{k_1+1, j_1, k_2+1, j_2}$ on the suffixes of these arcs.

Otherwise, exactly one of the leftmost events on \mathbf{R} and \mathbf{S} is an end point of a bond, and we have two symmetric cases ($QIs^{(1)}$ and $QIs^{(2)}$), one where the interaction arc is in \mathbf{R} , and the other where it is in \mathbf{S} . In the first case (middle diagram in Figure 7), let k_1 be the event in \mathbf{R} such that $i_1 \bullet k_1$ is an interaction arc, and $[i_2, k_2]$ is the longest subsequence in \mathbf{S} that $i_1 \bullet k_1$ subsumes, and k_2 is an event. The suffix of this substructure recurs on QI . We derive $QIs^{(1)}$ later.

To derive QIe_{i_1, j_1, i_2, j_2} , note that removing the arcs $i_1 \bullet j_1$ and $i_2 \bullet j_2$ yields the general case of $QI_{i_1+1, j_1-1, i_2+1, j_2-1}$ for the inner-section with an additional constraint that there has to be at least one bond in that region because the assumption is that the extracted arcs were interaction arcs. We can fulfill this constraint by excluding all cases where no bonds exist (i.e., considering only the two rightmost substructures of Figure 6).

To derive $QIs^{(1)}$, let k_1 be the leftmost event in the subsequence $[i_1+1, j_1-1]$. Note that such a k_1 is guaranteed to exist because first, $i_1 \bullet j_1$ subsumes $[i_2, j_2]$ and we know that i_2 is an event, i.e., the end point of either a bond (subsumed by $i_1 \bullet j_1$) or of an interaction arc. Then (see Figure 8) we define a new substructure, $QIaux^{(1)}$, after removing $i_1 \bullet j_1$ and the prefix of \mathbf{R} up to k_1 .

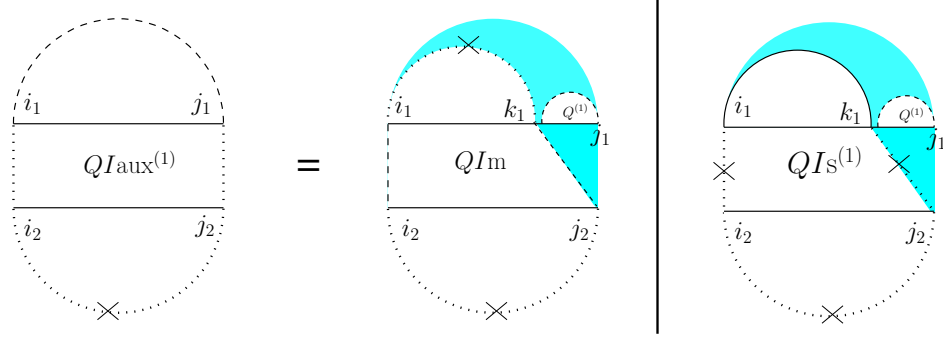


Figure 9: Two cases must be considered for the $QI_{\text{aux}}^{(1)}$ structure, in which the two end points of the bottom strand are events. For the top strand, only the leftmost end point is required to be an event. It can either be the end point of an arc (rightmost case) or not (leftmost case).

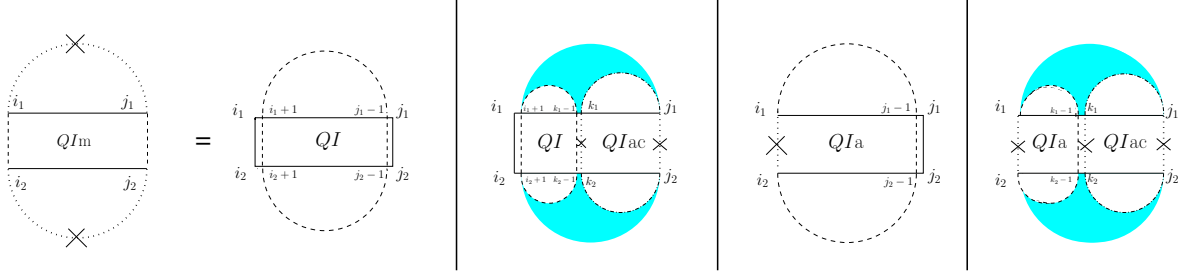


Figure 10: For computing QIm , since we know the four end points are events, but none of the two end points in one strand can form an arc, we must consider the four different cases shown above. For convenience, arcs of $QIac$ structure are shown with dash-dotted lines because it represents the sum of three structures in which each of the arcs could be present or not (we could replace the second and fourth cases with three cases, one for each term of Equation (16)).

To derive $QI_{\text{aux}}^{(1)}_{i_1, j_1, i_2, j_2}$, note that the context of its definition implies that i_1, i_2 and j_2 are all three events. Let, as shown in Figure 9, k_1 be the *last* event on $[i_1, j_1]$. Now, if $i_1 \bullet k_1$, then recur on $QIs^{(1)}$. Otherwise, k_1 is an event that does not pair with i_1 . We define a new substructure, QIm , where all four corners are events, and neither $i_1 \bullet j_1$ nor $i_2 \bullet j_2$ is allowed.

For computing QIm_{i_1, j_1, i_2, j_2} , since there are four corners each of which can be the end point of either a bond or of an arc, there might be at most sixteen possibilities. Upon combining some of those sixteen possibilities, we have to consider four mutually exclusive cases (see Figure 10). The first one is the case where $i_1 \circ i_2$ and $j_1 \circ j_2$ and the remaining part will be $QI_{i_1+1, j_1-1, i_2+1, j_2-1}$. That case corresponds to all four corner events being the end points of bonds. Since we assume there are no crossing bonds, we must have $i_1 \circ i_2$ and $j_1 \circ j_2$. In the second case, i_1 and i_2 are the end points of a bond, i.e., $i_1 \circ i_2$, but either j_1 or j_2 (or both) does not form a bond. That captures three out of the sixteen total possibilities. Since j_1 and j_2 are both events but do not form a bond, we define a term $QIac$ which is the sum of QIe and the two symmetric QIs 's, since they preserve the constraints that arise in the first term in the definition of QIa (see Figure 7). The prefix of this substructure is a general recursion on QI on the subsequences $[i_1 + 1, k_1 - 1]$ and $[i_2 + 1, k_2 - 1]$. The third case is the symmetric case of the second case, i.e., there is no bond between i_1 and i_2 , but $j_1 \circ j_2$. The prefix of this bond is a recursion on QIa . That captures three out of the sixteen total possibilities. Finally, the fourth case corresponds to either i_1 or i_2 (or both) does not form a bond and either j_1 or j_2 (or both) does not form a bond. That captures the remaining nine out of the sixteen total possibilities.

Putting all those together, we obtain

$$\text{QIa}_{i_1, j_1, i_2, j_2} = \sum_{k_1=i_1}^{j_1} \sum_{k_2=i_2}^{j_2} \text{QIac}_{i_1, k_1, i_2, k_2} \times \text{QI}_{k_1+1, j_1, k_2+1, j_2}, \quad (15)$$

$$\text{QIac}_{i_1, j_1, i_2, j_2} = \text{QIs}_{i_1, j_1, i_2, j_2}^{(1)} + \text{QIs}_{i_1, j_1, i_2, j_2}^{(2)} + \text{QIe}_{i_1, j_1, i_2, j_2}, \quad (16)$$

$$\text{QIe}_{i_1, j_1, i_2, j_2} = \begin{cases} 0 & j_1 - i_1 < 4 \\ & \text{or } j_2 - i_2 < 4 \\ M_{i_1, j_1, i_2, j_2} & \text{otherwise,} \end{cases} \quad (17)$$

$$M_{i_1, j_1, i_2, j_2} = \left(\text{QI}_{i_1+1, j_1-1, i_2+1, j_2-1} - \text{Q}_{i_1+1, j_1-1}^{(1)} \times \text{Q}_{i_2+1, j_2-1}^{(2)} \right) \times e^{\text{score}(i_1, j_1) + \text{score}(i_2, j_2)}, \quad (18)$$

$$\text{QIs}_{i_1, j_1, i_2, j_2}^{(1)} = \begin{cases} 0 & j_1 - i_1 < 4 \text{ or } j_2 < i_2 \\ \sum_{k_1=i_1+1}^{j_1-1} \text{Q}_{i_1+1, k_1-1}^{(1)} \times e^{\text{score}(i_1, j_1)} \times \text{QIaux}_{k_1, j_1-1, i_2, j_2}^{(1)} & \text{otherwise,} \end{cases} \quad (19)$$

$$\text{QIs}_{i_1, j_1, i_2, j_2}^{(2)} = \begin{cases} 0 & j_1 < i_1 \text{ or } j_2 - i_2 < 4 \\ \sum_{k_2=i_2+1}^{j_2-1} \text{Q}_{i_2+1, k_2-1}^{(2)} \times e^{\text{score}(i_2, j_2)} \times \text{QIaux}_{i_1, j_1, k_2, j_2-1}^{(2)} & \text{otherwise,} \end{cases} \quad (20)$$

$$\text{QIaux}_{i_1, j_1, i_2, j_2}^{(1)} = \sum_{k_1=i_1}^{j_1} \left(\text{QIs}_{i_1, k_1, i_2, j_2}^{(1)} + \text{QIm}_{i_1, k_1, i_2, j_2} \right) \times \text{Q}_{k_1+1, j_1}^{(1)}, \quad (21)$$

$$\text{QIaux}_{i_1, j_1, i_2, j_2}^{(2)} = \sum_{k_2=i_2}^{j_2} \left(\text{QIs}_{i_1, j_1, i_2, k_2}^{(2)} + \text{QIm}_{i_1, j_1, i_2, k_2} \right) \times \text{Q}_{k_2+1, j_2}^{(2)}, \quad (22)$$

$$\text{QIm}_{i_1, j_1, i_2, j_2} = \begin{cases} e^{\text{score}(i_1, i_2)} & i_1 = j_1 \text{ and } i_2 = j_2 \\ N_{i_1, j_1, i_2, j_2} & i_1 < j_1 \text{ and } i_2 < j_2 \\ 0 & \text{otherwise,} \end{cases} \quad (23)$$

$$N_{i_1, j_1, i_2, j_2} =$$

$$\begin{aligned} & e^{\text{score}(i_1, i_2) + \text{score}(j_1, j_2)} \times \text{QI}_{i_1+1, j_1-1, i_2+1, j_2-1} + \\ & e^{\text{score}(i_1, i_2)} \times \sum_{k_1=i_1+1}^{j_1} \sum_{k_2=i_2+1}^{j_2} \text{QI}_{i_1+1, k_1-1, i_2+1, k_2-1} \times \text{QIac}_{k_1, j_1, k_2, j_2} + \\ & e^{\text{score}(j_1, j_2)} \times \text{QIa}_{i_1, j_1-1, i_2, j_2-1} + \\ & \sum_{k_1=i_1}^{j_1} \sum_{k_2=i_2}^{j_2} \text{QIa}_{i_1, k_1, i_2, k_2} \times \text{QIac}_{k_1+1, j_1, k_2+1, j_2}. \end{aligned} \quad (24)$$

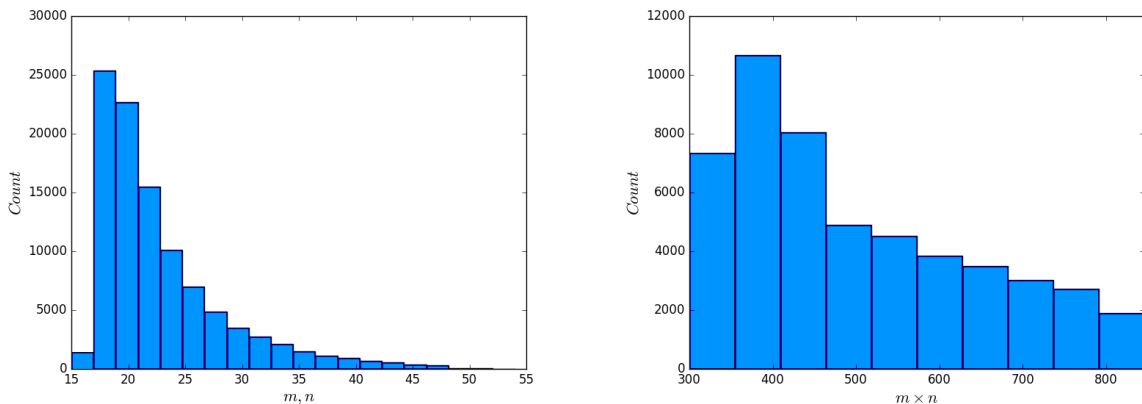


Figure 11: Distribution of the lengths (left) and product of the lengths (right) of all the RNA subsequences used in our experiment.

3 Results

To investigate the correlation between the scores of **BPPart** and **BPMax**, and those of **piRNA**, we used the RISE database [1] which combines information about interacting RNAs from multiple experiments. For the human dataset, we extracted all the interaction windows for those pairs that have this data in RISE. We eliminated the ones that contained a window with length less than 15 because they are too short for an unbiased comparison. Then, we sorted the remaining pairs based on the product of the lengths of the interacting windows. Finally, the first 50,500 pairs of sequences were chosen as our primary dataset for different experiments and analysis. Figure 11 shows the distribution of the sequence lengths in this dataset, and also the product of the lengths of the RNA subsequences in each pair.

We first ran **piRNA** on our primary dataset at 8 different temperatures, 37, 25, 13, 0, -40 , -80 , -130 , and -180 degrees Celsius. We also ran **BPPart** and **BPMax** on the dataset with different weights, i.e., c_i 's and c_i' 's. In general, we want to use the stack energies of the Turner model as a starting point for tuning the weights. Since the parameters form a projective space (invariant results with respect to scaling), we considered a fixed weight of 3 for **CG** (and **GC**). Using the experimentally computed stack energies of the Turner model, minimum and maximum values for the weights of **AU** and **GU** were computed. That is, to compute the maximum weight of **AU** (and **UA**), we consider the maximum released energy when **AU** (or **UA**) is stacked with another pair; this happens when **UA** is stacked with **CG** and 2.4 kcal/mol energy is released. Then, we considered the minimum value of released energy in an stack for **CG** or **GC** (for which we assumed a constant weight of 3), which is 1.4 kcal/mol . We derived the maximum weight of **AU** and **UA** as 5.143 by multiplying 2.4 by $\frac{3}{1.4}$. Finally, we made sure that the range of values that we explore for the weight of **AU** and **UA** contains this maximum value (we chose 5.5 as the upper-bound). For finding the minimum weight of **AU** and **UA**, we consider their minimum stack energy, which is 0.6 kcal/mol . Given the maximum energy of **CG**, namely 3.4 kcal/mol , the value of interest is computed as $0.6 \times \frac{3}{3.4} = 0.529$. However, for the sake of comprehensiveness and exploring the shape of the plots, we used much smaller lower-bounds— -4.5 and -3 —for **BPPart** and **BPMax**, respectively.

For all the combinations of weights of **AU** and **GU**, in steps of 0.5, we computed the Pearson and Spearman's rank correlations with the scores from **piRNA** at different temperatures. When computing the correlations, to normalize the scores from all algorithms, we divide them by the sum of the lengths of corresponding sequences, $L_R + L_S$. This normalization mitigates the effect of length bias on the computed correlations. This step is necessary because, generally, as the length of the pair of sequences increases the scores of all three algorithms increases, and if unnormalized scores are used, a biased higher correlation will be derived. Note that for partition functions, **piRNA** and **BPPart**, we used the \log of the scores; that is why we factor out the sum of the lengths for normalization. If the original values were to be used, we would have to take the

$(L_R + L_S)$ th roots of the scores. Figures 12 and 13 show the Pearson correlations for different combinations of weights of AU and GU at $-180^\circ C$ and $37^\circ C$. Figure 14 shows the scatter plots of the scores of **BPPart** and **piRNA** at these temperatures. In these plots, the red line shows the regression line that is fitted to the points by minimizing the mean squared error (MSE).

The optimum values of correlation for each temperature are presented in Tables 1 and 2. There is a high correlation between **piRNA** and **BPPart** as well as between **piRNA** and **BPMax**, especially when the temperature decreases which is due to a decrease in the role of thermodynamic entropy at the lower temperatures. Also, we computed the Pearson and Spearman’s rank correlation between **BPPart** and **BPMax** with their optimum weights at $37^\circ C$, which yielded values of 0.971 and 0.968, respectively.

Table 1: Pearson correlation between **piRNA** and **BPPart** and between **piRNA** and **BPMax** at different temperatures ($T^\circ C$).

Method \ T	37	25	13	0	-40	-80	-130	-180
BPPart	0.855	0.862	0.869	0.877	0.896	0.908	0.916	0.920
BPMax	0.836	0.846	0.855	0.864	0.884	0.895	0.901	0.904

Table 2: Spearman’s rank correlation between **piRNA** and **BPPart** and between **piRNA** and **BPMax** at different temperatures ($T^\circ C$).

Method \ T	37	25	13	0	-40	-80	-130	-180
BPPart	0.864	0.867	0.871	0.876	0.889	0.896	0.901	0.901
BPMax	0.830	0.835	0.841	0.847	0.862	0.871	0.877	0.877

To make sure that the optimization results are not data dependent, we conducted the same experiments for randomly generated sequences. To factor out the effect of length, for each pair in our primary dataset, we generated a pair of random sequences with the same lengths as those of the pair in our primary dataset. The shape of the plots are very similar to those for the primary dataset, for both **BPPart** and **BPMax** (see Figures 15, 16). For **BPPart**, the optimum weights at $-180^\circ C$ are the same (1.0, 1.0, and 3 for AU, GU, and CG, respectively) and at $37^\circ C$ the optimum weights we had earlier (0.5, 1.0, and 3) are ranked 3rd for the random dataset with only 0.003 difference in the Pearson correlation from the optimum, achieved by weights of (1.0, 1.0, and 3) for the random dataset. Similarly, for **BPMax**, the optimum values for the two datasets are the same at $-180^\circ C$ (1.0, 2.0, and 3), and at $37^\circ C$, the optimum values of weights for the primary dataset (1.0, 1.5, and 3) are ranked 2nd for the random dataset with only 0.008 difference in the Pearson correlation from the optimum, achieved by weights of (1.0, 2.0, and 3) for this dataset.

Although the shape of the plots and the peaks were almost the same for the primary dataset and the random dataset, the best correlations for the random one were considerably less than those of the primary one. Table 3 shows the Pearson and the Spearman’s rank correlation of **BPPart** and **BPMax** with **piRNA** at $-180^\circ C$ and $37^\circ C$ for the random dataset.

Table 3: Pearson and Spearman’s rank correlation between **piRNA** and **BPPart** and between **piRNA** and **BPMax** at $-180^\circ C$ and $37^\circ C$ for the random input sequences.

Method \ T	Pearson		Spearman	
	$37^\circ C$	$-180^\circ C$	$37^\circ C$	$-180^\circ C$
BPPart	0.761	0.825	0.753	0.801
BPMax	0.716	0.785	0.702	0.751

Finally to better understand the behavior of the surface around the higher values in the correlation plots of Figures 12 and 13, we computed the Shannon entropy for the values above a threshold. Figure 17 shows

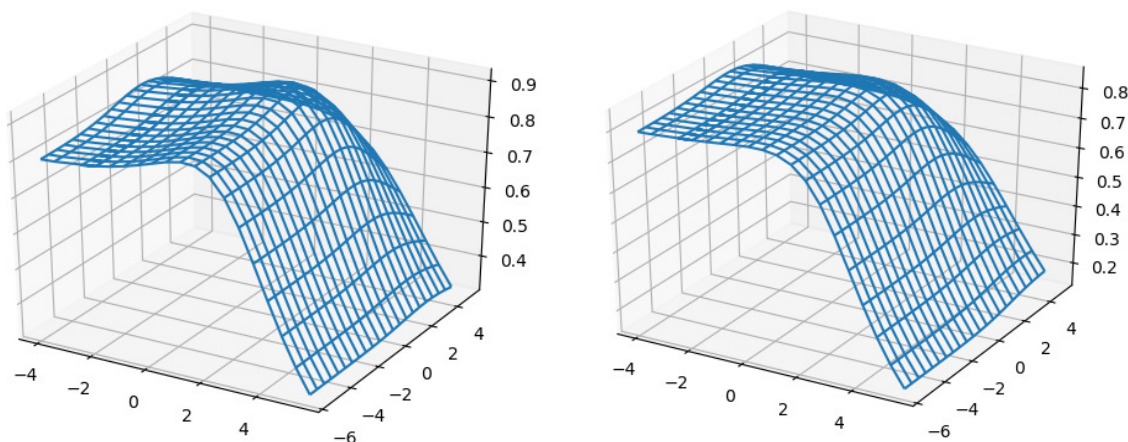


Figure 12: Pearson correlation between **piRNA** and **BPPart** (vertical axis), on the primary dataset, at -180°C (left) and 37°C (right) for different values of constant factors (weights) for **AU** (left axis) and **GU** (right axis). The weight of **CG** pair is fixed at 3.

these values for the top 30 values of Pearson and Spearman’s rank correlation at each temperature.

Discussion

The Gibbs free energy

$$\Delta G = \Delta H - T\Delta S \quad (25)$$

is composed of a term ΔH called enthalpy that does not depend on temperature and a term $T\Delta S$ called entropy that linearly depends on temperature T . Intuitively, enthalpy is the chemical energy that is often released upon formation of chemical bonds such as base pairing. Entropy, on the other hand, captures the size of all possible spatial conformations for a fixed secondary structure. In other words, entropy captures the amount of 3D freedom of the molecule. A base pair brings enthalpy down, hence favorable from enthalpy point of view, and decreases freedom (entropy), hence unfavorable from entropy point of view. These two opposing objectives are combined linearly through the temperature coefficient.

In the full thermodynamic model, we consider both terms. In the base pair counting, we consider only a simplistic enthalpy term. Partition function for the full thermodynamic model is

$$\sum_{s \in \mathcal{S}^I} e^{-\Delta G/RT}, \quad (26)$$

in which R is the gas constant. Note that

$$-\frac{\Delta G}{T} = -\frac{\Delta H}{T} + \Delta S, \quad (27)$$

and as $T \rightarrow 0$, $-\Delta H/T \rightarrow \infty$ and the contribution of ΔS is diminished to 0 since it is finite. Hence, at low temperatures, the effect of entropy becomes negligible, and we expect to see strong correlation between the base pair counting model and full thermodynamic model.

Figure 18 shows the Pearson correlations between **BPPart** and **BPMax** scores with **piRNA** scores for a fixated combination of weights that results in the highest correlation at 37°C . For **BPPart** the chosen weights are 0.5, 1.0, and 3 for **AU**, **GU**, and **CG**, respectively, while the corresponding weights for **BPMax** are 1.0, 1.5, and 3.

Perfectly conforming with the theory, we see higher correlations at low temperatures. That somewhat validates our implementations as **piRNA** was written totally independently about 10 years ago. Moreover, as

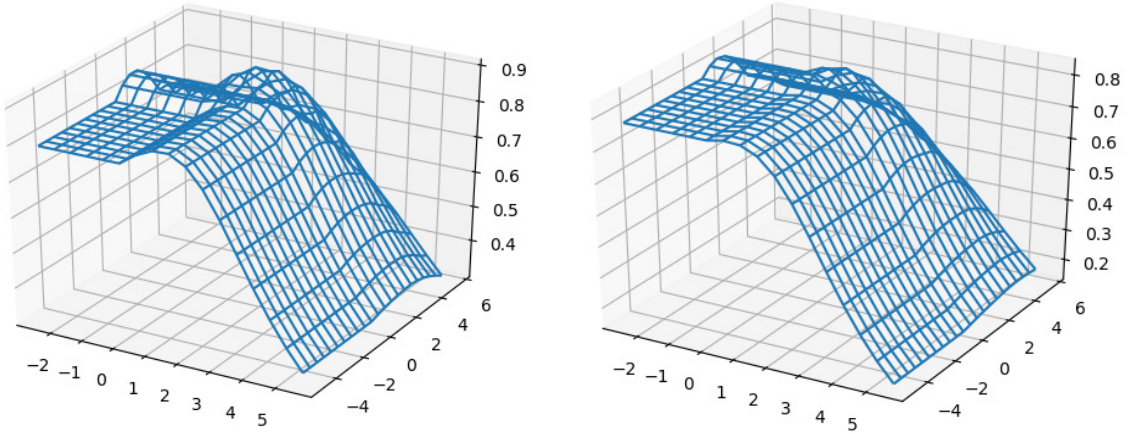


Figure 13: Pearson correlation between **piRNA** and **BpMax** (vertical axis), on the primary dataset, at -180°C (left) and 37°C (right) for different values of constant factors (weights) for **AU** (left axis) and **GU** (right axis). The weight of **CG** pair is fixed at 3.

can be seen in Figures 12 and 13, the surface around the optimum value for higher temperatures becomes flatter. Figure 17, which shows the entropy of the top 30 correlation values, confirms this observation. This means the correlation values are less sensitive to a change in the weights of the base pairs as the temperature increases; this conforms with the theory because at higher temperatures, the thermodynamic entropy increases and the total score of **piRNA** becomes less sensitive to the energy released by pairings. It is worth mentioning that having less Shannon entropy for the top values at higher temperatures decreases the possibility of having universal optimum values for the weights of the base pairs.

Another noticeable characteristic of the plots in Figures 12 and 13 is the region in which the scores of both **AU** and **GU** are non-positive. This region for **BpMax** is flat because when both of these pairs are penalized (or not rewarded when their score is zero), the algorithm simply avoids making such pairs because it is trying to maximize the score. Therefore, it only tries to maximize the number of **CG** pairs, which is independent of the scores (penalty in this case) of the other two types of base pairs. This also applies to the case where one of the base pairs has a non-positive score; in that case, **BpMax** works independently of the score of that base pair. So, as soon as any of the scores becomes non-positive, **BpMax** remains constant along the corresponding axis. For **BpPart**, however, the story is different because it simply counts all the possible pairings and even if the score of a base pair becomes negative, it does not ignore counting that.

Moreover, **BpPart** has a higher correlation than **BpMax** does, which comes with the price of a $6\times$ increase in empirical running time. Also, as Figure 17 shows, the Shannon entropy for the top 30 values is less in **BpMax** and the gap between them grows as temperature decreases; this shows that **BpPart** has a flatter region around the optimum value and its optimum value is less sensitive to changes in the weights. Meanwhile, having a steeper surface in **BpMax** which has less entropy, increases the possibility of having more stable and universal optimum values for the weights. As mentioned earlier, the running time difference between the two is noticeable: **BpMax** is about $6\times$ faster than **BpPart**. Hence, we now have three choices in increasing order of computational cost: **BpMax**, **BpPart**, and **piRNA**. The computation time increases by about $6\times$ and $225\times$, respectively, from one to the next.

Finally, based on the results of the experiments on both the primary dataset and the random one, we see that although the shapes of the optimization plots and the optimum weights are very similar, the correlation values are much less for the random dataset. This observation suggests that probably the interaction regions are more complementary than the random sequences of the same size because in these regions the effect of the energy released by pairing probably becomes more significant than the energy added by an increase in the entropy on the final score of **piRNA**. That could explain why we get a higher correlation in such regions with **BpPart** and **BpMax**, which mainly depend on the base pairs. This hypothesis has to be thoroughly investigated in the future.

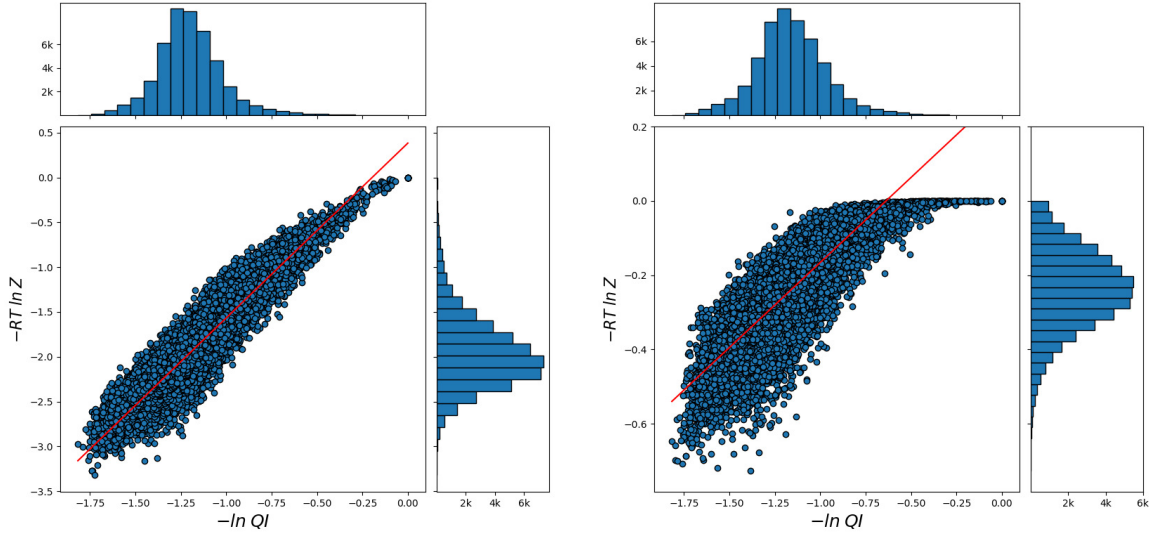


Figure 14: Scatter plots of the values of the corresponding axis for each sample (interaction windows of a pair of RNAs from RISE dataset) at -180°C (left) and 37°C (right). In both plots, the red line is a straight regression line fitted to the points by minimizing mean square error (MSE).

Application of BPPart in Human Biology

One of the use-cases of **BPPart** and **BPMax**, among others, is making predictions about the consequences of a slight change in the RNA sequences. This information becomes helpful for various domains and tasks, such as synthetic biology and studying the mutations. Between **BPMax** and **BPPart**, the latter is much more sensitive to small changes in the sequence, because it considers all possible structures that the two interacting sequences might form. Therefore, even a missense mutation might make a tangible difference in the computed **BPPart** score.

To verify this hypothesis, we used **BPPart** to study the effects of known missense mutations, provided by Ensembl, in the interaction regions of some RISE pairs. Given a pair of interacting RNAs in RISE for which the information about the interacting regions is provided, we retrieved the data of all the reported missense mutations of those regions from Ensembl API. Also, we got the phenotypic consequences of each mutation from Ensembl. Finally, we computed the **BPPart** score for the original sequence of one of the interacting regions and each of the mutated versions of the other sequence. Among all the generated scores for a pair, we found the outliers using the interquartile range. These outliers, represent a mutation in one of the interacting RNAs, which falls within the interacting region, that causes a great difference in the interaction score. In the rest of this section, we almost-randomly pick and narrate two of such cases that we observed, among many discovered ones. In Appendix A, we report 65 such pairs that have been discovered using this pipeline after analyzing more than one million pairs of sequences that have been generated after applying the known missense mutations to over 15,200 pairs of interacting sequences reported in the RISE database. Further study of each of these pairs and more comprehensive study of effect of nonsense mutations on RRI would be a next step in the future.

Traces of TRAF3 in CADASIL

Cerebral Autosomal Dominant Arteriopathy with Subcortical Infarcts and Leukoencephalopathy (CADASIL) is an inherited condition in which the muscle cells of small blood vessels, especially the ones in the brain, gradually die and cause many impairments, such as stroke, cognitive impairment, and mood disorders in the elderly [48]. It has been shown that mutation in NOTCH3, which resides on the reverse strand of chromosome 19, is responsible for this condition in people with this genetic disorder [49]. NOTCH3 and

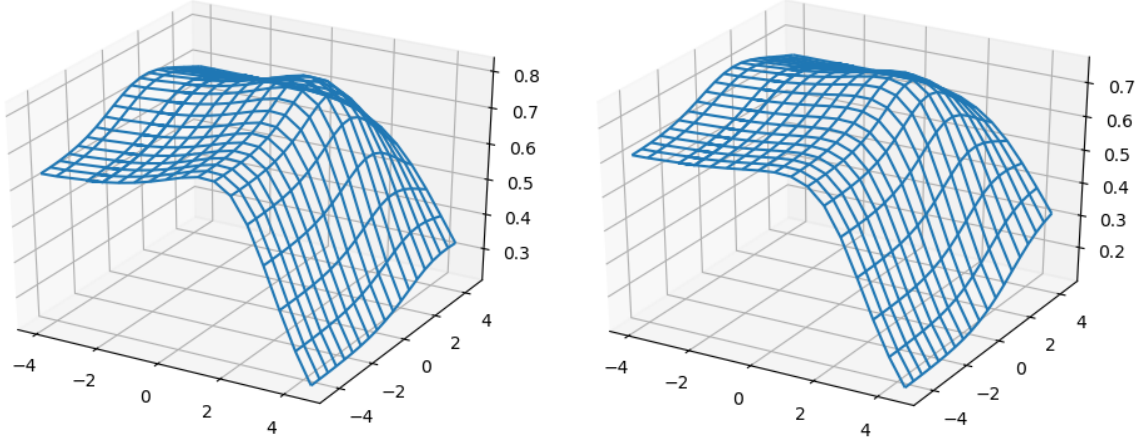


Figure 15: Pearson correlation between **piRNA** and **BPPart** (vertical axis), on the randomly generated dataset, at $-180^{\circ}C$ (left) and $37^{\circ}C$ (right) for different values of constant factors (weights) for *AU* (left axis) and *GU* (right axis). The weight of *CG* pair is fixed at 3.

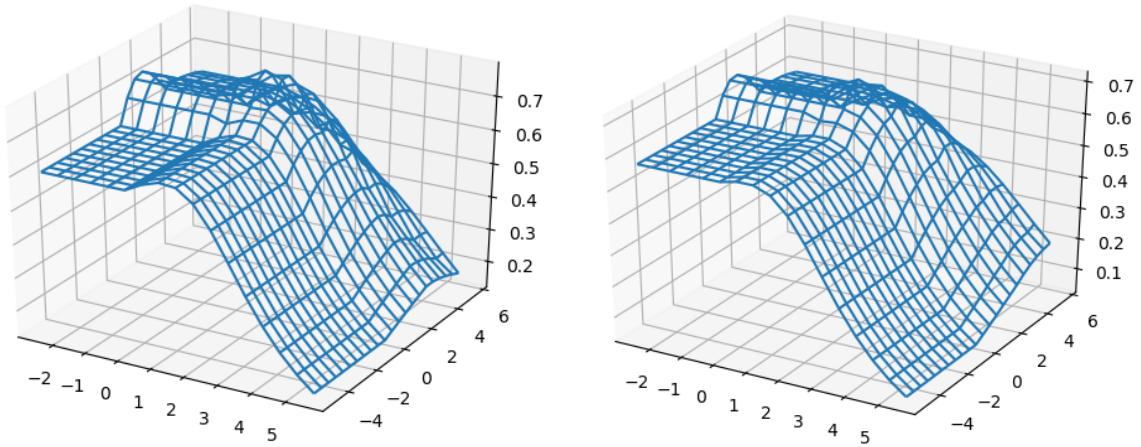


Figure 16: Pearson correlation between **piRNA** and **BPPart** (vertical axis), on the randomly generated dataset, at $-180^{\circ}C$ (left) and $37^{\circ}C$ (right) for different values of constant factors (weights) for *AU* (left axis) and *GU* (right axis). The weight of *CG* pair is fixed at 3.

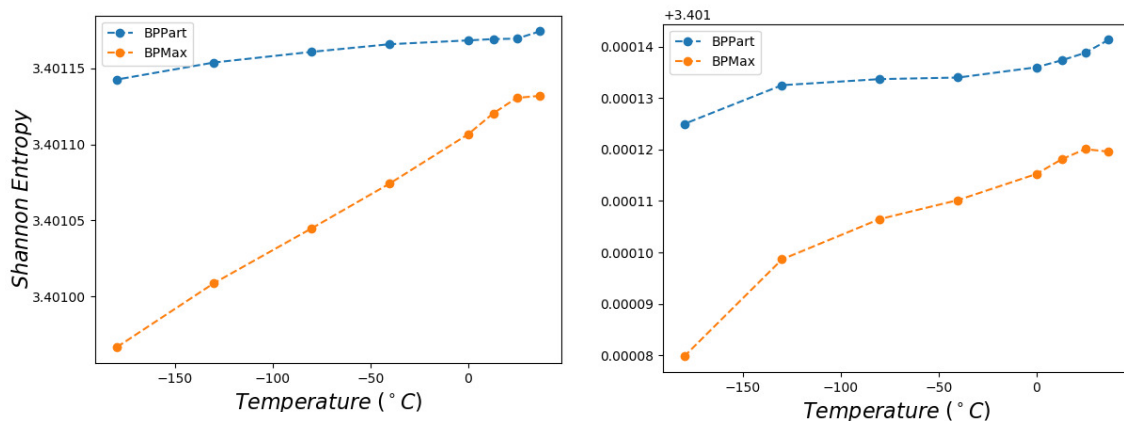


Figure 17: Shannon entropy for the top 30 Pearson (left) and Spearman's rank (right) correlation values at different temperatures for BPPart and BPMMax.

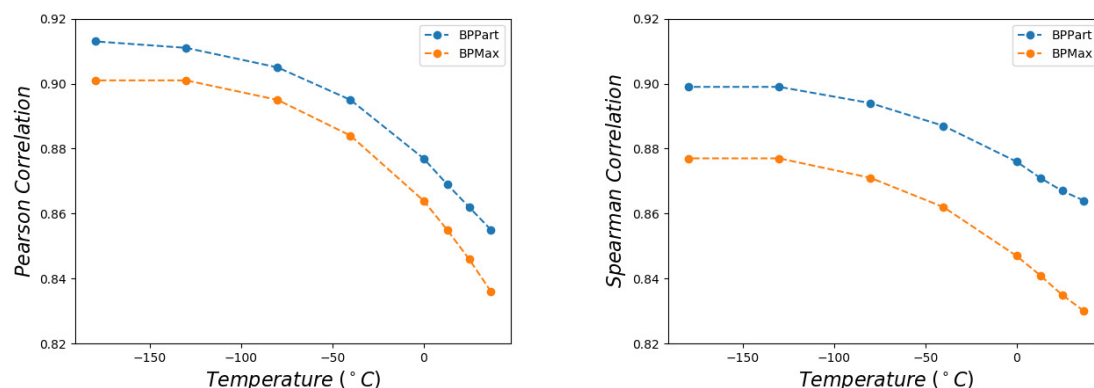


Figure 18: Pearson correlation (left) and Spearman's rank correlation (right) between piRNA and BPPart and between piRNA and BPMMax at different temperatures.

TRAF3 are a pair of interacting RNAs that have been reported in RISE. One of the missense mutations in NOTCH3 that has been reported to be contributing to CADASIL [50] lies within the interacting region of this gene, from loci 15,161,520 to 15,161,543 (according to GRCh38 assembly of human genome), with TRAF3. Interestingly, this mutation, which replaces nucleotide *C* with *G* at loci 15,161,526 of chromosome 19, causes a dramatic increase in the score of BPPart such that it makes it an outlier when the aforementioned procedure is followed. TRAF3 is a gene that has been reported to play a role in angiogenesis [51, 52]. A noticeable increase in the score of BPPart increases the chance that these two RNAs interact and cause post-transcriptional conditions that affect the translation rate of TRAF3 which possibly contributes to the phenotypic consequences of CADASIL. Further evaluation and verification of this hypothesis requires further experimental analysis.

Traces of SNORD3D in Parkinson's Disease

SNORD3D is a small nucleolar RNA which has been detected not long ago [53] with which no specific task or annotation is associated in the literature yet. According to the RISE database, one of the genes that interacts with this snoRNA is GBA which resides on the reverse strand of chromosome 1. Mutations in

GBA has been reported to play a role in Parkinson’s disease which is a brain disorder that affects movement and often causes tremors. One of the GBA mutations that is reported to be linked with Parkinson’s disease lies within the interaction region of this gene, from loci 155,239,966 to 155,239,984 (according to GRCh38 assembly of human genome), with SNORD3D. This specific mutation of GBA, which changes the nucleotide *G* to *C* at loci 155,239,972 of chromosome 1, is one of the cases that is detected as an outlier using our aforementioned analysis of **BPPart** scores. This mutation, when applied to GBA, decreases its score of interaction with SNORD3D, which might cause the interaction to occur much less than the normal case. This possibly leads to a change in the expression of GBA protein. According to KEGG, GBA is a member of Other glycan degradation, Sphingolipid metabolism, Metabolic pathways, and Lysosome pathways [54]. Therefore, we hypothesize the role of SNORD3D in some or all of those pathways, particularly, the ones that are closely related to Parkinson’s disease. Further evaluation of this hypothesis requires further experimental data and analysis.

4 Conclusion

We revisited the problems of partition function and structure prediction for interacting RNAs. We simplified the energy model and instead considered only simple weighted base pair counting to obtain **BPPart** for the partition function and **BPMax** for structure prediction. As a result, **BPPart** runs about 225 \times and **BPMax** runs about 1300 \times faster than **piRNA** does. Hence, we gained significant speedup by potentially sacrificing accuracy.

To evaluate practical accuracy of both new algorithms, we computed the Pearson and rank correlations at different temperatures between the results of **BPPart**, **BPMax**, and **piRNA** on 50,500 experimentally characterized RRI in the RISE database [1]. **BPPart** and **BPMax** results correlate well with those of **piRNA** at low temperatures. At the room and body temperatures, there is considerable correlation and therefore, significant information in the results of **BPPart** and **BPMax**.

We conclude that both **BPPart** and **BPMax** capture a significant portion of the thermodynamic information. Both tools can be used as filtering steps in more sophisticated RRI prediction pipelines. Also, the information captured by **BPPart** and **BPMax** can possibly be complemented with machine learning techniques in the future for more accurate predictions. We now have three choices for RRI thermodynamics in increasing computational cost: **BPMax**, **BPPart**, and **piRNA**. Depending on the application and the trade-off between time and accuracy, one can be chosen.

Finally, we show that **BPPart** might be useful to explain how some mutations lead to some specific phenotypic consequences. We presented two new hypotheses about the roles of TRAF3 in CADASIL and SNORD3D in lipid processing pathways and/or Parkinson’s disease.

4.1 Conflict of interest statement.

None declared.

References

- [1] J. Gong, D. Shao, K. Xu, Z. Lu, Z. J. Lu, Y. T. Yang, and Q. C. Zhang. RISE: a database of RNA interactome from sequencing experiments. *Nucleic Acids Res.*, Oct 2017.
- [2] Hamidreza Chitsaz, Raheleh Salari, S.Cenk Sahinalp, and Rolf Backofen. A partition function algorithm for interacting nucleic acid strands. *Bioinformatics*, 25(12):i365–i373, 2009. Also ISMB/ECCB proceedings.
- [3] D.H. Mathews, J. Sabina, M. Zuker, and D.H. Turner. Expanded sequence dependence of thermodynamic parameters improves prediction of RNA secondary structure. *J. Mol. Biol.*, 288:911–940, May 1999.
- [4] S. Cao and S.J. Chen. Predicting RNA pseudoknot folding thermodynamics. *Nucleic Acids Res.*, 34:2634–2652, 2006.
- [5] Robert M. Dirks and Niles A. Pierce. A partition function algorithm for nucleic acid secondary structure including pseudoknots. *Journal of Computational Chemistry*, 24(13):1664–1677, 2003.
- [6] R. Nussinov, G. Piecznik, J. R. Grigg, and D. J. Kleitman. Algorithms for loop matchings. *SIAM Journal on Applied Mathematics*, 35:68–82, 1978.
- [7] M. S. Waterman and T. F. Smith. RNA secondary structure: A complete mathematical analysis. *Math. Biosc.*, 42:257–266, 1978.
- [8] Michael Zuker and Patrick Stiegler. Optimal computer folding of large RNA sequences using thermodynamics and auxiliary information. *Nucleic Acids Research*, 9(1):133–148, 1981.
- [9] E. Rivas and S.R. Eddy. A dynamic programming algorithm for RNA structure prediction including pseudoknots. *J. Mol. Biol.*, 285(5):2053–2068, 1999.
- [10] J.S. McCaskill. The equilibrium partition function and base pair binding probabilities for RNA secondary structure. *Biopolymers*, 29:1105–1119, 1990.
- [11] Anne Wenzel, Erdiç Akbaşlı, and Jan Gorodkin. Rsearch: fast rna–rna interaction search using a simplified nearest-neighbor energy model. *Bioinformatics*, 28(21):2738–2746, 2012.
- [12] M. Andronescu, Z.C. Zhang, and A. Condon. Secondary structure prediction of interacting RNA molecules. *J. Mol. Biol.*, 345:987–1001, Feb 2005.
- [13] S.H. Bernhart, H. Tafer, U. Mückstein, C. Flamm, P.F. Stadler, and I.L. Hofacker. Partition function and base pairing probabilities of RNA heterodimers. *Algorithms Mol Biol*, 1:3, 2006.
- [14] Robert M. Dirks, Justin S. Bois, Joseph M. Schaeffer, Erik Winfree, and Niles A. Pierce. Thermodynamic analysis of interacting nucleic acid strands. *SIAM Review*, 49(1):65–88, 2007.
- [15] M. Rehmsmeier, P. Steffen, M. Hochsmann, and R. Giegerich. Fast and effective prediction of microRNA/target duplexes. *RNA*, 10:1507–1517, Oct 2004.
- [16] Roumen A. Dimitrov and Michael Zuker. Prediction of hybridization and melting for double-stranded nucleic acids. *Biophysical Journal*, 87:215–226, 2004.
- [17] N.R. Markham and M. Zuker. UNAFold: software for nucleic acid folding and hybridization. *Methods Mol. Biol.*, 453:3–31, 2008.
- [18] D.D. Pervouchine. IRIS: intermolecular RNA interaction search. *Genome Inform*, 15:92–101, 2004.
- [19] Can Alkan, Emre Karakoc, Joseph H. Nadeau, S. Cenk Sahinalp, and Kaizhong Zhang. RNA-RNA interaction prediction and antisense RNA target search. *Journal of Computational Biology*, 13(2):267–282, 2006.

- [20] Ronny Lorenz, Stephan H Bernhart, Christian Hoener Zu Siederdisen, Hakim Tafer, Christoph Flamm, Peter F Stadler, and Ivo L Hofacker. Viennarna package 2.0. *Algorithms for Molecular Biology*, 6(1):26, 2011.
- [21] Laura DiChiacchio, Michael F Sloma, and David H Mathews. Accessfold: predicting rna–rna interactions with consideration for competing self-structure. *Bioinformatics*, 32(7):1033–1039, 2015.
- [22] Yuki Kato, Tatsuya Akutsu, and Hiroyuki Seki. A grammatical approach to RNA-RNA interaction prediction. *Pattern Recognition*, 42(4):531–538, 2009.
- [23] Fenix W. D. Huang, Jing Qin, Christian M. Reidys, and Peter F. Stadler. Partition function and base pairing probabilities for RNA-RNA interaction prediction. *Bioinformatics*, 25(20):2646–2654, 2009.
- [24] U. Mückstein, H. Tafer, J. Hackermüller, S.H. Bernhart, P.F. Stadler, and I.L. Hofacker. Thermodynamics of RNA-RNA binding. *Bioinformatics*, 22:1177–1182, May 2006.
- [25] S.P. Walton, G.N. Stephanopoulos, M.L. Yarmush, and C.M. Roth. Thermodynamic and kinetic characterization of antisense oligodeoxynucleotide binding to a structured mRNA. *Biophys. J.*, 82:366–377, Jan 2002.
- [26] Anke Busch, Andreas S. Richter, and Rolf Backofen. IntaRNA: Efficient prediction of bacterial sRNA targets incorporating target site accessibility and seed regions. *Bioinformatics*, 24(24):2849–2856, 2008.
- [27] Hamidreza Chitsaz, Rolf Backofen, and S.Cenk Sahinalp. biRNA: Fast RNA-RNA binding sites prediction. In S.L. Salzberg and T. Warnow, editors, *Workshop on Algorithms in Bioinformatics (WABI)*, volume 5724 of *LNBI*, pages 25–36, Berlin, Heidelberg, 2009. Springer-Verlag.
- [28] Azra Krek, Dominic Grün, Matthew N Poy, Rachel Wolf, Lauren Rosenberg, Eric J Epstein, Philip MacMenamin, Isabelle Da Piedade, Kristin C Gunsalus, Markus Stoffel, et al. Combinatorial microrna target predictions. *Nature genetics*, 37(5):495, 2005.
- [29] Michael Kertesz, Nicola Iovino, Ulrich Unnerstall, Ulrike Gaul, and Eran Segal. The role of site accessibility in microrna target recognition. *Nature genetics*, 39(10):1278, 2007.
- [30] Jan Krüger and Marc Rehmsmeier. Rnahybrid: microrna target prediction easy, fast and flexible. *Nucleic acids research*, 34(suppl_2):W451–W454, 2006.
- [31] Yuanji Zhang. miru: an automated plant mirna target prediction server. *Nucleic acids research*, 33(suppl_2):W701–W704, 2005.
- [32] Ivo L Hofacker, Walter Fontana, Peter F Stadler, L Sebastian Bonhoeffer, Manfred Tacker, and Peter Schuster. Fast folding and comparison of rna secondary structures. *Monatshefte für Chemie/Chemical Monthly*, 125(2):167–188, 1994.
- [33] Jin-Wu Nam, Olivia S Rissland, David Koppstein, Cei Abreu-Goodger, Calvin H Jan, Vikram Agarwal, Muhammed A Yildirim, Antony Rodriguez, and David P Bartel. Global analyses of the effect of different cellular contexts on microrna targeting. *Molecular cell*, 53(6):1031–1043, 2014.
- [34] Doron Betel, Anjali Koppal, Phaedra Agius, Chris Sander, and Christina Leslie. Comprehensive modeling of microrna targets predicts functional non-conserved and non-canonical sites. *Genome biology*, 11(8):R90, 2010.
- [35] Martin Reczko, Manolis Maragkakis, Panagiotis Alexiou, Ivo Grosse, and Artemis G Hatzigeorgiou. Functional microrna targets in protein coding sequences. *Bioinformatics*, 28(6):771–776, 2012.
- [36] Dimos Gaidatzis, Erik van Nimwegen, Jean Hausser, and Mihaela Zavolan. Inference of mirna targets using evolutionary conservation and pathway analysis. *BMC bioinformatics*, 8(1):69, 2007.
- [37] Ángela Riffo-Campos, Ismael Riquelme, and Priscilla Brebi-Mieville. Tools for sequence-based mirna target prediction: what to choose? *International journal of molecular sciences*, 17(12):1987, 2016.

- [38] Kevin C Miranda, Tien Huynh, Yvonne Tay, Yen-Sin Ang, Wai-Leong Tam, Andrew M Thomson, Bing Lim, and Isidore Rigoutsos. A pattern-based method for the identification of microRNA binding sites and their corresponding heteroduplexes. *Cell*, 126(6):1203–1217, 2006.
- [39] Justin Bo-Kai Hsu, Chih-Min Chiu, Sheng-Da Hsu, Wei-Yun Huang, Chia-Hung Chien, Tzong-Yi Lee, and Hsien-Da Huang. mirtar: an integrated system for identifying mirna-target interactions in human. *BMC bioinformatics*, 12(1):300, 2011.
- [40] Wenlong Xu, Anthony San Lucas, Zixing Wang, and Yin Liu. Identifying microRNA targets in different gene regions. *BMC bioinformatics*, 15(7):S4, 2014.
- [41] Anke Busch, Andreas S Richter, and Rolf Backofen. Intarna: efficient prediction of bacterial srna targets incorporating target site accessibility and seed regions. *Bioinformatics*, 24(24):2849–2856, 2008.
- [42] Martin Mann, Patrick R Wright, and Rolf Backofen. Intarna 2.0: enhanced and customizable prediction of rna–rna interactions. *Nucleic acids research*, 45(W1):W435–W439, 2017.
- [43] Stephanie Kehr, Sebastian Bartschat, Peter F Stadler, and Hakim Tafer. Plexy: efficient target prediction for box c/d snornas. *Bioinformatics*, 27(2):279–280, 2010.
- [44] Hakim Tafer, Stephanie Kehr, Jana Hertel, Ivo L Hofacker, and Peter F Stadler. Rnasnoop: efficient target prediction for h/aca snornas. *Bioinformatics*, 26(5):610–616, 2010.
- [45] Brian Tjaden. Targetrna: a tool for predicting targets of small rna action in bacteria. *Nucleic acids research*, 36(suppl_2):W109–W113, 2008.
- [46] Sinan Uğur Umu and Paul P Gardner. A comprehensive benchmark of rna–rna interaction prediction tools for all domains of life. *Bioinformatics*, 33(7):988–996, 2017.
- [47] R. Nussinov and A. Jacobson. Fast algorithm for predicting the secondary structure of single stranded RNA. *Proc. Nat. Acad. Sci. USA*, 77(11):6309–6313, 1980.
- [48] Ilaria Di Donato, Silvia Bianchi, Nicola De Stefano, Martin Dichgans, Maria Teresa Dotti, Marco Duerling, Eric Jouvent, Amos D Korczyn, Saskia AJ Lesnik-Oberstein, Alessandro Malandrini, et al. Cerebral Autosomal Dominant Arteriopathy with Subcortical Infarcts and Leukoencephalopathy (CADASIL) as a model of small vessel disease: update on clinical, diagnostic, and management aspects. *BMC medicine*, 15(1):41, 2017.
- [49] Anne Joutel, Christophe Corpechot, Anne Ducros, Katayoun Vahedi, Hugues Chabriat, Philippe Mouton, Sonia Alamowitch, Valérie Domenga, Michaelle Cecillion, Emmanuelle Marechal, et al. Notch3 mutations in cerebral autosomal dominant arteriopathy with subcortical infarcts and leukoencephalopathy (CADASIL), a mendelian condition causing stroke and vascular dementia. *Annals of the New York Academy of Sciences*, 826(1):213–217, 1997.
- [50] Anne Joutel, Christophe Corpechot, Anne Ducros, Katayoun Vahedi, Hugues Chabriat, Philippe Mouton, Sonia Alamowitch, Valérie Domenga, Michaelle Cécillion, Emmanuelle Maréchal, et al. Notch3 mutations in CADASIL, a hereditary adult-onset condition causing stroke and dementia. *Nature*, 383(6602):707–710, 1996.
- [51] Hong Ming Hu, Karen O’Rourke, Mark S Boguski, and Vishua M Dixit. A novel RING finger protein interacts with the cytoplasmic domain of CD40. *Journal of Biological Chemistry*, 269(48):30069–30072, 1994.
- [52] Almin I Lalani, Carissa R Moore, Chang Luo, Benjamin Z Kreider, Yan Liu, Herbert C Morse, and Ping Xie. Myeloid cell TRAF3 regulates immune responses and inhibits inflammation and tumor development in mice. *The Journal of Immunology*, 194(1):334–348, 2015.
- [53] Mitchell Guttman, Ido Amit, Manuel Garber, Courtney French, Michael F Lin, David Feldser, Maite Huarte, Or Zuk, Bryce W Carey, John P Cassady, et al. Chromatin signature reveals over a thousand highly conserved large non-coding RNAs in mammals. *Nature*, 458(7235):223–227, 2009.

- [54] S. Tsuji, P. V. Choudary, B. M. Martin, S. Winfield, J. A. Barranger, and E. I. Ginns. Nucleotide sequence of cDNA containing the complete coding sequence for human lysosomal glucocerebrosidase. *J. Biol. Chem.*, 261(1):50–53, Jan 1986.

Appendix A: Disease-causing mutations that incur an outlier BPPart score with a known RRI partner

disease_causing_outliers												
rise_id	id1	id2	name1	name2	mutachr	strand	position	orig_new	phenotype	IQR	deviation	
0	51273	ENSG00000136068	ENSG00000148734	FLNB	NPFFR1	1 chr3	+	58124465	C G	FLNB-Related Spectrum Disorders	1.28964999999999	--0.07925000000000892
1	50415	ENSG00000140470	ENSG00000123607	ADAMTS17	TTC21B	2 chr2	-	165912577	G A	Jeune thoracic dystrophy;Joubert Syndrome	0.894625000000005	--0.40788749999999396
2	61127	ENSG00000142197	ENSG00000105983	DOPEY2	LMBR1	2 chr7	-	156762153	A C	Triphalangeal thumb polysyndactyly syndrome	1.16395	0.257249999999999
3	58037	ENSG00000136937	ENSG00000156873	NCBP1	PHKG2	2 chr16	+	30756397	T G	GLYCOGEN STORAGE DISEASE IXc;GLYCOGEN	0.286300000000004	0.872049999999994
4	62929	ENSG00000189056	ENSG00000240771	RELN	ARHGEF25	1 chr7	-	103472745	G A	Lissencephaly, Recessive	0.342799999999997	--2.3163000000000054
5	63228	ENSG00000117174	ENSG00000151929	ZNHIT6	BAG3	2 chr10	+	119672416	C G	"Primary familial hypertrophic cardiomyopathy;Myofit	0.578600000000002	0.749600000000001
6	49884	ENSG00000183943	ENSG00000156873	PRKX	PHKG2	2 chr16	+	30756397	T G	GLYCOGEN STORAGE DISEASE IXc;GLYCOGEN	0.659300000000002	0.744249999999994
7	57967	ENSG00000158813	ENSG00000130520	EDA	LSM4	1 chrX	+	69957098	G T	Hypohidrotic X-linked ectodermal dysplasia;Hyp	1.0635	0.170150000000007
8	54237	ENSG00000187942	ENSG00000101463	LDLRAD2	SYNDIG1	1 chr1	+	21822492	C G	Schwartz Jampel syndrome type 1;Dyssegmental Dy	0.525700000000001	2.01465
9	61924	ENSG00000234465	ENSG00000156873	PINLYP	PHKG2	2 chr16	+	30756397	T G	GLYCOGEN STORAGE DISEASE IXc;GLYCOGEN	0.400399999999998	1.273100000000001
10	53477	ENSG00000184937	ENSG00000270641	WT1	TSIX	1 chr11	-	32434971	A G	WILMS TUMOR, ANIRIDIA, GENITOURINARY ANO	0.660525000000007	0.515787499999973
11	53142	ENSG00000177628	ENSG00000281000	GBA	SNORD3D	1 chr1	-	155239972	G C	Parkinson disease, late-onset;Parkinson disease, lat	0.616700000000002	--0.5511749999999935
12	53382	ENSG00000128645	ENSG00000118972	HOXD1	FGF23	2 chr12	-	4370564	C A	HYPOPHOSPHATEMIC RICKETS, AUTOSOMAL D	0.513374999999996	--0.1644124999999974
13	58401	ENSG00000188493	ENSG00000187535	C19orf54	IFT140	2 chr16	-	1520311	G A	Renal dysplasia, retinal pigmentary dystrophy, cereb	0.192399999999992	--1.63745000000000154
14	51135	ENSG00000161594	ENSG00000125730	KLHL10	C3	2 chr19	-	6709754	G A	HEMOLYTIC UREMIC SYNDROME, ATYPICAL, SU	0.292400000000001	--0.07389999999999475
15	54836	ENSG00000196712	ENSG00000163626	NF1	COX18	1 chr17	+	31225205	G A	Neurofibromatosis, type 1	0.667774999999999	--0.31838750000000715
16	51249	ENSG00000163513	ENSG00000171444	TGFB2	MCC	1 chr3	+	30606953	G T	Thoracic aortic aneurysm and aortic dissection	1.158175	--0.15768750000000153
17	50639	ENSG00000147027	ENSG00000164190	TMEM47	NIPBL	2 chr5	+	36876831	G C	Cornelia de Lange syndrome	1.123324999999999	--1.5334125000000043
18	65472	ENSG00000240490	ENSG00000177098	RN7SL277P	SCN4B	2 chr11	-	118135693	G A	Long QT syndrome;Romano-Ward syndrome	0.259449999999994	--1.4180250000000072
19	52715	ENSG00000177663	ENSG00000196943	IL17RA	NOP9	1 chr22	+	17102166	G A	Familial Candidiasis, Recessive	0.998850000000004	--0.09147499999998843
20	53657	ENSG00000186684	ENSG00000156873	CYP27C1	PHKG2	2 chr16	+	30756397	T G	GLYCOGEN STORAGE DISEASE IXc;GLYCOGEN	0.411999999999999	1.0973
21	62243	ENSG00000197430	ENSG00000179915	OPALIN	NRXN1	2 chr2	-	51028694	A C	Pitt-Hopkins-like syndrome	0.123999999999995	--0.53400000000000131
22	55306	ENSG00000174469	ENSG00000112319	CNTNAP2	EYA4	1 chr7	+	148415896	C A	Pitt-Hopkins-like syndrome;CORTICAL DYSPLAS	0.457799999999999	--0.6374000000000066
23	60574	ENSG00000166813	ENSG00000147144	KIF7	CCDC120	1 chr15	-	89642233	G C	Acrocallosal syndrome, Schinzel type"	1.113225000000001	0.218512499999989
24	54812	ENSG00000054654	ENSG00000239900	SYNE2	ADSL	2 chr22	+	40364966	G A	ADENYLOSUCCINASE DEFICIENCY;Adenylosuccii	1.1359	--0.3874499999999941
25	49578	ENSG00000138095	ENSG00000156521	LRPPRC	TYSND1	1 chr2	-	43887016	G A	Leigh syndrome	0.099999999999994	--1.67090000000000103
26	50039	ENSG00000074181	ENSG00000131323	NOTCH3	TRAF3	1 chr19	-	15161526	C G	Cerebral autosomal dominant arteriopathy with subc	0.729074999999995	0.8911875000000015
27	57248	ENSG00000151929	ENSG00000279659	BAG3	RP11-177G23	1 chr10	+	119672297	C G	Myofibrillar myopathy, BAG3-related;Dilated Cardio	0.253600000000006	1.496024999999999
28	50475	ENSG00000141576	ENSG00000166147	RNF157	FBN1	2 chr15	-	48432910	G A	MARFAN SYNDROME	0.906750000000002	--0.240524999999991
29	50754	ENSG00000124155	ENSG00000164588	PIGT	HCN1	2 chr5	-	45262033	G C	Early infantile epileptic encephalopathy	0.498525000000001	-0.038812500000006
30	62856	ENSG00000130234	ENSG00000169604	ACE2	ANTXR1	2 chr2	+	69152194	G A	HEMANGIOMA, CAPILLARY INFANTILE, SUSCEP	0.279900000000005	-1.355049999999999
31	61544	ENSG00000221869	ENSG00000166813	CEBPD	KIF7	2 chr15	-	89648308	A G	Acrocallosal syndrome, Schinzel type"	0.550025000000005	0.614062499999989
32	54646	ENSG00000127914	ENSG00000229807	AKAP9	XIST	1 chr7	+	92002229	G A	Long QT syndrome	0.934375000000003	-0.019387499999994
33	48374	ENSG00000100345	ENSG00000178996	MYH9	SNX18	1 chr22	-	36387950	C G	MYH9-related disorder;Nonsyndromic Hearing Loss,	0.480525	0.739812499999999
34	57452	ENSG00000101974	ENSG00000156709	ATP11C	AIFM1	2 chrX	-	130131756	G A	Deafness, X-linked 5"	0.576374999999999	-0.539487500000007
35	64673	ENSG00000075624	ENSG00000180182	ACTB	MED14	1 chr7	-	5529594	G A	BARAITSER-WINTER SYNDROME 1	0.971800000000002	-0.159500000000008
36	64673	ENSG00000075624	ENSG00000180182	ACTB	MED14	1 chr7	-	5529624	A G	BARAITSER-WINTER SYNDROME 1;BARAITSER-	0.971800000000002	0.399899999999988
37	64673	ENSG00000075624	ENSG00000180182	ACTB	MED14	1 chr7	-	5529624	A C	BARAITSER-WINTER SYNDROME 1;BARAITSER-	0.971800000000002	0.141399999999999
38	51417	ENSG00000164619	ENSG00000105576	BMPER	TNPO2	1 chr7	+	33905094	G A	Diaphanospondylodysostosis	0.296549999999996	-0.3862000000000017
39	62328	ENSG00000232316	ENSG00000170289	RP1-124C6	CNGB3	2 chr8	-	86574166	G A	Stargardt Disease, Recessive;Achromatopsia	0.725475000000001	-0.24951249999998
40	53469	ENSG00000183196	ENSG00000166233	CHST6	ARIH1	1 chr16	-	75477044	T G	Macular corneal dystrophy Type I	0.450999999999993	0.446200000000005
41	61935	ENSG00000183873	ENSG00000108306	SCN5A	FBXL20	1 chr3	-	38550198	A G	Paroxysmal familial ventricular fibrillation;Roman	0.237174999999993	1.6237375
42	59210	ENSG00000207110	ENSG00000173575	RNU1-106P	CHD2	2 chr15	+	93020077	G A	Epileptic encephalopathy, childhood-onset;Epileptic	0.112300000000005	-0.015799999999992
43	64224	ENSG00000165283	ENSG00000183873	STOML2	SCN5A	2 chr3	-	38633255	G A	Congenital long QT syndrome	0.395350000000008	-0.459074999999984
44	48691	ENSG00000110092	ENSG00000076248	CCND1	UNG	2 chr12	+	109110470	C G	Immunodeficiency with Hyper-IgM	0.374000000000002	1.551599999999999
45	62535	ENSG00000105576	ENSG00000110799	TNPO2	VWF	2 chr12	-	5981833	G T	von Willebrand disorder	0.867199999999997	-0.5235000000000013
46	62535	ENSG00000105576	ENSG00000110799	TNPO2	VWF	2 chr12	-	5981833	G A	von Willebrand disorder	0.867199999999997	-1.133600000000001
47	51170	ENSG00000149136	ENSG00000100234	SSRP1	TIMP3	2 chr22	+	32862388	G A	Pseudoinflammatory fundus dystrophy	0.221199999999996	-1.612100000000001
48	58057	ENSG00000185920	ENSG00000162341	PTCH1	TPCN2	1 chr9	-	95469067	A G	Gorlin syndrome	0.7308000000000016	1.063449999999997
49	58057	ENSG00000185920	ENSG00000162341	PTCH1	TPCN2	1 chr9	-	95469088	G A	Gorlin syndrome	0.7308000000000016	-0.740349999999964
50	53161	ENSG00000183864	ENSG00000075624	TOB2	ACTB	2 chr7	-	5529624	A C	BARAITSER-WINTER SYNDROME 1;BARAITSER-	0.956874999999997	0.639812500000005
51	64674	ENSG00000075624	ENSG00000251497	ACTB	RP11-197N18	1 chr7	-	5529594	G A	BARAITSER-WINTER SYNDROME 1	0.705799999999996	-0.336849999999998

disease_causing_outliers													
52	53476	ENSG00000184937	ENSG00000252680	WT1	RNA5SP449	1 chr11	-	32434971	A	G	WILMS TUMOR, ANIRIDIA, GENITOURINARY ANO	0.821150000000003	0.399924999999982
53	53476	ENSG00000184937	ENSG00000252680	WT1	RNA5SP449	1 chr11	-	32434980	C	G	Drash syndrome;WILMS TUMOR, ANIRIDIA, GEN	0.821150000000003	0.681824999999989
54	65331	ENSG00000163686	ENSG00000075624	ABHD6	ACTB	2 chr7	-	5527786	G	A	BARAITSER-WINTER SYNDROME 1	0.394350000000003	-0.161524999999998
55	51773	ENSG00000164588	ENSG00000105576	HCN1	TNPO2	1 chr5	-	45695893	T	C	Early infantile epileptic encephalopathy	0.374624999999998	0.822012500000042
56	53048	ENSG00000158467	ENSG00000105610	AHCYL2	KLF1	2 chr19	-	12885877	G	A	Congenital dyserythropoietic anemia	0.503799999999998	-0.714799999999997
57	61065	ENSG00000249158	ENSG00000130164	PCDHA11	LDLR	2 chr19	+	11132534	C	G	Familial hypercholesterolemia	0.647699999999986	0.641475000000028
58	53557	ENSG00000184634	ENSG00000140323	MED12	DISP2	1 chrX	+	71137822	G	T	FG syndrome	0.388400000000004	-1.47449999999998
59	58063	ENSG00000072501	ENSG00000202324	SMC1A	RNA5SP366	1 chrX	-	53376265	G	A	Cornelia de Lange syndrome	0.564299999999989	-0.531150000000011
60	56554	ENSG00000201861	ENSG00000136068	RNA5SP298	FLNB	2 chr3	+	58077256	G	A	FLNB-Related Disorders	0.698800000000006	-1.1193
61	56554	ENSG00000201861	ENSG00000136068	RNA5SP298	FLNB	2 chr3	+	58077266	T	G	BOOMERANG DYSPLASIA;BOOMERANG DYSPLA/	0.698800000000006	0.633399999999998
62	56554	ENSG00000201861	ENSG00000136068	RNA5SP298	FLNB	2 chr3	+	58077271	G	A	ATELOSTEOGENESIS, TYPE III;Atelosteogenesis t	0.698800000000006	-0.962399999999988
63	53823	ENSG00000124107	ENSG00000169071	SLPI	ROR2	2 chr9	-	91733195	G	A	Brachydactyly;Robinow syndrome	0.744950000000003	-0.730199999999982
64	49653	ENSG00000101577	ENSG00000198931	LPIN2	APRT	1 chr18	-	2920313	G	A	MAJEED SYNDROME	0.404200000000003	-1.661575