

RNA Folding: Local Versus Global Optimization

Max Ward (20748588)

*This report is submitted as partial fulfilment
of the requirements for the Honours Programme of the
School of Computer Science and Software Engineering,
The University of Western Australia,
2014*

Abstract

I am definitely going to need to write this at some point. This is a short report on how to use the `cshonours.cls` class to prepare dissertations using the latest \LaTeX version, $\text{\LaTeX}2\text{e}$. This class is based on the standard class `report.cls`.

Keywords: Honours, report, dissertation, UWA, RNA, bioinformatics

CR Categories: Not, really, sure

Acknowledgements

Going to need something here too. This class is designed to produce reports that look the same as those produced by the older `cshonours.sty` style for \LaTeX 2.09, which was modified by Nick Spadaccini from a style provided by Ken Wessen.

Contents

Abstract	ii
Acknowledgements	iii
1 Introduction	1
1.1 DNA and RNA	1
1.2 Dynamic Programming Techniques	4
1.2.1 Fundamental Algorithms	4
1.2.2 Pseudoknots	9
1.3 Accuracy	10
1.4 Alternative Techniques	11
1.4.1 Soft Computing	11
1.4.2 Context Free Grammars	12
1.5 Locally Optimal Structure Prediction	13
1.6 State of the Art: Global Optimization	14
2 Methology	16
2.1 Materials	16
2.1.1 Environment	16
2.1.2 Software	16
2.1.3 Testing Set	17
2.2 Merging RNA Intervals	17
2.2.1 Top-Down Selection	18
2.2.2 Bottom-Up Selection	19
2.2.3 MFE Selection	19
2.2.4 Weighted Activity Selection	19
2.2.5 Comparison of RNA Interval Selection Algorithms	20

2.2.6	RNA Prediction With Windows	20
-------	---------------------------------------	----

List of Tables

List of Figures

1.1	The structure and composition of DNA. Diagram taken from "Essential Cell Biology" [1].	2
1.2	RNA transcription. Diagram taken from "Essential Cell Biology" [1].	3
1.3	RNA secondary structure as described in the Nussinov algorithm. Taken from the original publication [21].	5
1.4	Diagram of faces used in the Zuker algorithm. Taken from original publication [39].	7

CHAPTER 1

Introduction

1.1 DNA and RNA

Deoxyribonucleic Acid (DNA) is the basic genetic building block upon which the classification of genetic material into genes and chromosomes is based. The role of DNA as the hereditary unit of genetics was determined in the 1940s [1]. Soon thereafter, Watson & Crick [35] published a highly acclaimed paper describing the fundamental chemical structure of DNA. In it, they outlined a double helix formation which has since become as iconic as it is canonical (see Figure 1.1). Each strand of the helix Watson & Crick discovered is essentially a chain of ‘nucleotides’ which are made of a sugar-phosphate backbone, attached to a single ‘base’. The bases of each strand form hydrogen bonds which hold the double helix together. The most astonishing and important of their findings was that these bases bond in a reciprocal fashion. They described four bases: Adenine (A), which always bonds to Thymine (T), and Guanine (G), which always bonds to Cytosine (C).

The reciprocal bonding relationships between bases is what allows replication to occur; a copy of the DNA can be made by simply allowing the correct bases to bond to one of the strands making up its helix. This gives a model for inheritance and cellular replication. However, there remains the question of how DNA can actually code for protein. Proteins are made up amino acids bonded in a specific sequence [1]. The DNA must therefore code for amino acids. This code, which can be thought of as the ‘digital’ representation for the ‘analogue’ protein used by our cells, needs to be carried to ribosomes which translate it into protein [1]. This is a task carried out by Ribonucleic Acid (RNA). RNA is very much like DNA in that it can bond reciprocally to another strand with matching bases. The main difference is that it is single stranded in structure, and has Uracil (U) in place of Thymine [1]. It is important to note that in RNA molecules G and U pairings are also possible. RNA bonds to DNA and, in a sense, reads it. This results in the production of a copy of the DNAs genetic payload. This

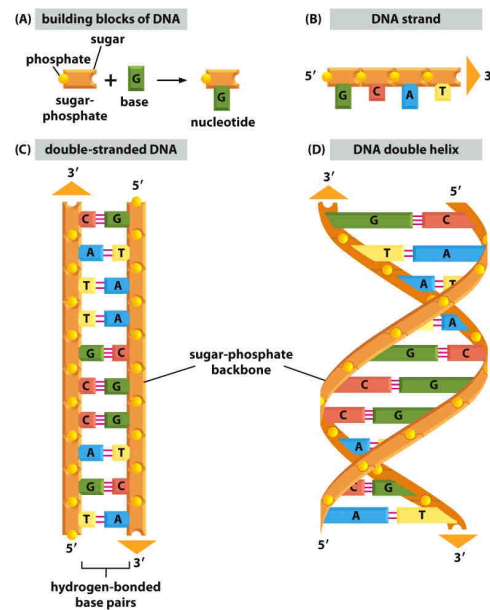


Figure 1.1: The structure and composition of DNA. Diagram taken from "Essential Cell Biology" [1].

‘downloaded’ information is then carried away to be translated into protein [1]. An example of this is depicted in Figure 1.2, in which we see a Messenger RNA molecule bonding to and thus making a copy of a section of DNA. As depicted in Figure 1.2, the 3’ end of a DNA or RNA molecule is the end onto which new nucleotides are added. The 5’ end is chemically stable, and nucleotides are not usually appended to it [1].

For many years the conventional wisdom was that DNA contained genes which coded for functional proteins used by the cell [1]. Though this is undoubtedly true, there was a problem: much of the human genome, and the genomes of other species, contains DNA which does not appear to code for anything [4]. Many theories have been put forward to explain this. It was argued that this ‘junk’ DNA is the perennial build-up of mutation, and that natural selection simply cannot act with strong enough selective force to cull this free-loading DNA [4]. Surprisingly, much of this non-coding DNA is actually transcribed into RNA, despite having no apparent function [16]. As it turns out, RNA is more than a simple messenger for encoded proteins. Recent research has found myriad important functions for RNA. For example, RNA can act as a catalyst for RNA splicing and peptide bond formation, and can also alter the regulation of genes [37]. It seems that much of our genome contains templates for non-coding RNAs

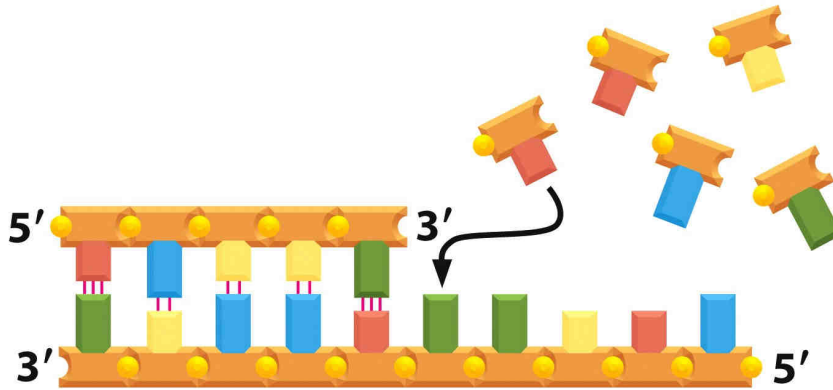


Figure 1.2: RNA transcription. Diagram taken from "Essential Cell Biology" [1].

(ncRNAs). These RNAs perform essential cellular functions without actually being translated into protein at any point in their life-cycle [16]. Because of its inherently single stranded nature, RNA forms bonds with itself, folding into secondary and tertiary structures [6].

It is axiomatic that chemical structure is tantamount to biological function; RNA is no exception. For this reason there has and continues to be an intense interest in predicting the secondary structure and tertiary structure of RNA molecules. This is in part because it will elucidate the underlying principles of RNA structure formation and function [6], but also because it will allow the detection and classification of unknown RNAs, enable prediction of novel RNA function, and assist the design of new RNA based drugs [5]. In fact, RNA is an extremely versatile molecule, and as such is attractive from both an engineering and computational point of view. Small combinatorial computation problems have been solved by representing the solution set using RNAs. Furthermore, a theory of computation has been put forward using self assembling RNA molecules [5]. As if to comment on the upheaval of a protein-centric view of biology in recent years, researchers have found that RNA is capable of supporting all the processes required for life without the need of protein [5]. The secondary structure of RNA is also highly conserved during evolution, indicating its importance [11]. Secondary and tertiary structures can be treated hierarchically, as a result it is possible to predict the secondary structure of an RNA without understanding the tertiary structure. The tertiary structure in turn builds upon the secondary structure [32]. This paper will focus on secondary structure prediction.

It holds to reason that an algorithm for RNA secondary structure prediction

can never be realised if we do not understand how these structures form, or their general morphology. For this reason it is important to understand how true RNA secondary structures can be determined, and the limitations of these techniques. DNA and RNA molecules can be analysed using X-ray crystallographic methods. These types of approaches work because the wavelengths of some X-rays are the same as the dimensions of DNA and RNA inter-atomic bonds. The diffraction of X-ray light by these molecules can thus be observed and their structures can subsequently be inferred by analysis of the resulting data. Nuclear Magnetic Resonance (NMR) is another technique which can be applied to the analysis of DNA/RNA. It relies on the spin of atoms when in a magnetic field. These spin signals can be used to determine the atomic composition and topology of a molecule. This has the advantage of not requiring the molecule under analysis to be crystallized before analysis. Arguably this gives a better *in vivo* view of RNAs/DNAs, which are fundamentally flexible structures. NMR also has some disadvantages; for instance, it is less accurate than X-ray crystallography, and cannot be used on extremely large molecules. The reason these techniques cannot be used for all RNA structural assays is that they are extremely expensive and time consuming [20].

RNA secondary structure prediction techniques can be broadly broken into two categories: those that use auxiliary information to assist in prediction, and those that predict structure *ex nihilo*—that is, with nothing but the ‘proband’ sequence we require a structure for. The former approach typically does consensus matching between some sequences for which a user already knows the secondary structures, and a sequence for which the structure is unknown [11]. In this paper I investigate the latter approach because it requires deeper knowledge about why and how RNAs fold. Also, it is the more general of the two.

1.2 Dynamic Programming Techniques

1.2.1 Fundamental Algorithms

The first such algorithms were based on relatively naive brute force. All possible secondary structures were enumerated and the one with the most bonds was selected as the solution [22]. While being very simplistic, these first approaches introduce an important assumption: RNA molecules will form energetically stable secondary structures. Maximising bonds is a crude but nonetheless accurate measure of energetic stability, as every bond increases the stability of a structure [22]. In the late 1970s, when the first large RNA molecules were being successfully sequenced, Nussinov et al. [22] introduced an algorithm based on loop matching

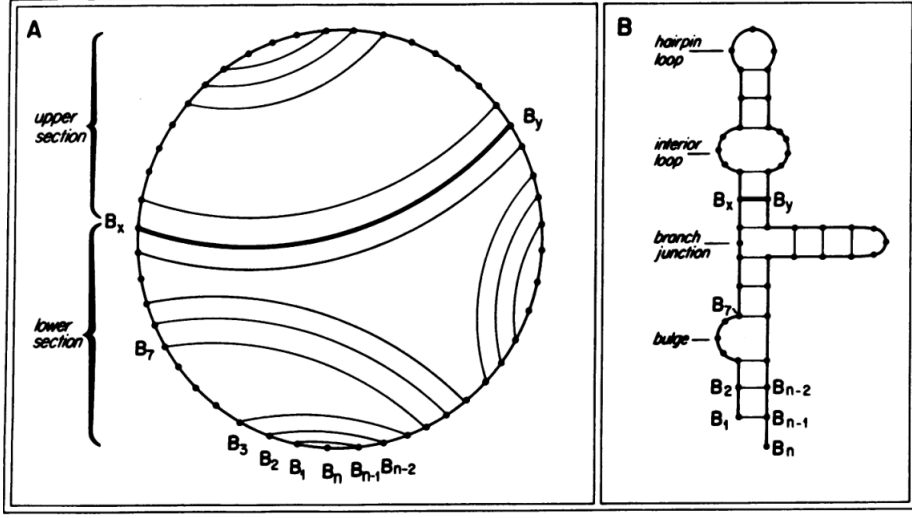


Figure 1.3: RNA secondary structure as described in the Nussinov algorithm. Taken from the original publication [21].

for bonding pairs. Their algorithm attempted to find a single structure with the maximal number of bonds using dynamic programming, with the restriction that all bonding pairs had to be entirely nested. It did this in $O(N^3)$ time and using $O(N^2)$ space. Thence Nussinov & Jacobson [21] introduced a refined version of the same algorithm and began testing it against experimentally verified RNA secondary structures. They had mixed success; Transfer RNAs (tRNAs) were conspicuous in their difference from predicted structures.

Because of its dynamic programming nature, this algorithm performs recursive decompositions of the RNA and builds larger structures out of repeated substructures. A natural representation of this is depicted in Figure 1.3. Part A of Figure 1.3 shows bonds as arcs across a circular graph. In it, we see the nested nature of the structures being explored by the Nussinov algorithm. Part B shows how these structures translate to actual RNAs, and how these appear in vivo. It also introduces the standard decompositions of secondary structures, namely the hairpin loop, the interior loop, and the branch junction or multi loop. Unlabelled in the diagram are stems; these are stacked base pairings, for example B_1, B_{n-1} and B_2, B_{n-2} .

$$M(i, j) = \max \{A, B, C, D\} \quad (1.1)$$

$$A = M(i, j - 1)$$

$$B = M(i + 1, j)$$

$$C = M(i + 1, j - 1) + W(i, j)$$

$$D = \max \{M(i, k) + M(k + 1, j)\} \text{ when } i < k < j$$

In the recurrence relation shown in Equation 1.1 the first two cases (*A* and *B*) find the score associated with not allowing *i* and *j* to bond. The case *C* conversely determines the score given that *i* and *j* are bonded. The final case *D* computes the score associated with a bifurcation. A bifurcation here means decomposition of the RNA into two separate structures between. This recurrence relation implies a $O(N^3)$ worst case time complexity and a $O(N^2)$ space complexity, as an $O(N^2)$ state space (all combinations of *i* and *j*) is explored with a linear time recurrence relation. In the original algorithm a constant $p = 3$ was introduced that indicated the minimum size of a hair-pin loop as real RNAs typically do not have hair-pin loops of fewer bases. The recurrence relation presented here has also been modified for the sake of clarity (cases *A* and *B* can be merged into case *D*) but the logic of the algorithm is equivalent.

This algorithm can also be extended to accommodate a more advanced energy model. Instead of weighting each bond equally they can be weighted according to the proportion they are expected to contribute to the molecule's stability [21]. When considering the value of a bond, it might be given greater weight if it adds to the formation of a stem (a stabilizing structure), or given lower weight if it forms an internal loop or bulge as these generally destabilize RNA molecules [21]. Unfortunately it is hard to find good values for such weights, and determining which substructure a bond contributes to requires backtracking in the modified algorithm presented by Nussinov & Jacobson.

The reader should note that the Nussinov algorithm is old technology and is no longer used for the prediction of RNA secondary structures. I have presented it in detail here because it forms the basis for the most widely used algorithm today, which I shall then discuss since it in turn is the basis for my own algorithms.

Soon after the work of Nussinov & Jacobson, Zuker & Stiegler [39] described an altered version of the same algorithm which, instead of maximising base pairs, minimized the free energy of the secondary structure. This was done by introducing a number of thermodynamic rules for canonical structures like hairpin loops, internal bulges, multiloops, unbonded base pairs, and stacked base pairs. The algorithm is similar to the Nussinov algorithm but adds another mutually recursive dynamic programming recurrence to inject a complex and relatively comprehensive energy system. The original energy system used is borrowed from the work of Studnicka et al. [30] who presented a complex algorithm which predicted similar RNA secondary structures, albeit with much worse asymptotic and implementation complexities.

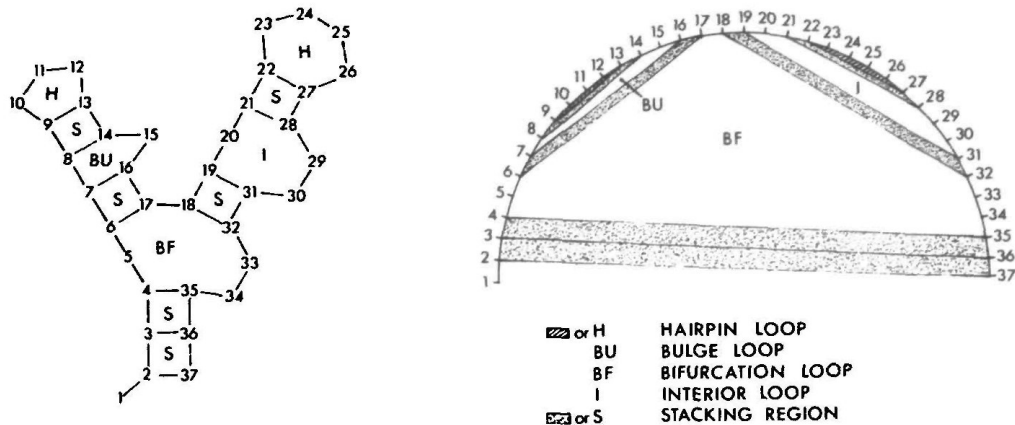


Figure 1.4: Diagram of faces used in the Zuker algorithm. Taken from original publication [39].

First I shall introduce some useful terminology which should clarify aspects of Zuker & Stieglers algorithm. The bases of an RNA molecule can be thought of as vertices in a graph. Edges between such vertices can be represented as chords on a semicircular diagram (Figures 1.3 and 1.4), such chords are not allowed to touch. A chord is admissible if the bonds connected by it are chemically valid bonds, and an admissible structure is a structure whose graph contains only admissible bonds. Thence one can define a face of such a graph as any planar region bounded on all sides; these faces represent the basic substructures of an RNA molecule [39]. The folding algorithm of Zuker & Stieglers considers such faces as the basic contributing factor to a molecules stability, unlike the algorithm of Nussinov & Jacobson which considers only individual bonds.

Let $E(F)$ represent the energy of a face F ; impossible structures are given an energy value of infinity, for example hairpin loops with less than three bases in the intervening loop region. In addition let $V(i, j)$ be defined as the minimum free energy of all structures in which bases i and j are bonded, and let $W(i, j)$ represent the minimum free energy of all structures contained within bases i and j inclusive. Note that for $W(i, j)$ there may or may not be a bond between bases i and j . Also if i and j cannot bond then $V(i, j) = \infty$. Finally note that $FH(i, j)$ represents a hairpin loop structure from i to j , and that $FL(i, j, i', j')$ is defined as the region bounded by the bonds i, j and i', j' . Examples of these decompositions are shown diagrammatically in the right half of Figure 1.4. The labelled regions show faces in a semicircular graph representing a strand of RNA.

In the accompanying left half of the figure, the same RNA structure is shown as it would appear in a real RNA rather than in a purely diagrammatic depiction.

$$\begin{aligned}
V(i, j) &= \min \{E1, E2, E3\} \\
E1 &= E(FH(i, j)) \\
E2 &= \min \{E(FL(i, j, i', j')) + V(i', j')\} \text{ where } i < i' < j' < j \\
E3 &= \min \{W(i + 1, i') + W(i' + 1, j1)\} \text{ where } i + 1 < i' < j - 2
\end{aligned} \tag{1.2}$$

As shown by the definition provided in Equation 1.2, $V(i, j)$ is computed by minimizing three cases. The first case considers the bond between i and j closing off a hairpin loop (H in Figure 1.4). The second accounts for cases in which i and j are bonded. This can result in a bulge (BU in Figure 1.4), internal loop (I in Figure 1.4), or the continuation of a stacking region (S in Figure 1.4) between the interior bond i', j' . The third and final case considers bifurcations (BF in Figure 1.4)

$$\begin{aligned}
W(i, j) &= \min \{W(i + 1, j), W(i, j1), V(i, j), E4\} \\
E4 &= \min \{W(i, i') + W(i' + 1, j)\} \text{ where } i < i' < j - 1
\end{aligned} \tag{1.3}$$

In Equation 1.3 we show the recurrence for $W(i, j)$ as described by Zuker & Stiegler. Again there are three cases. The first two cases $W(i + 1, j)$ and $W(i, j1)$ should be thought of as a single case which consider situations in which there is no bond between i and j . This is similar to cases *A* and *B* from the Nussinov algorithm (Equation 1.1). The final case considers taking the bond from i to j . This final case allows for bifurcations in which two bonding pairs split the structure into two sections. The final minimum free energy of the best structure is defined by $W(1, n)$, where n is the length of the RNA molecule. It should be noted that the free energy for small molecules (fewer than 6 nucleotides in length) can easily be precomputed, and forms the base case of the given recurrence relations. Because of its efficiency ($O(N^3)$ time and $O(N^2)$ space), robustness, and extensibility, this method, even today, is still the most popular technique available. The most widely used packages for RNA secondary structure prediction are all implementations of the Zuker algorithm [24, 17].

The Zuker algorithm suffers a major shortcoming, however. Because all bonding regions are assumed to be nested, it cannot handle the case of ‘pseudoknots’. Pseudoknots are structures in which a bonding pair may have its first base inside another bonding pair, and the other base outside said bonding pair. In short, it is not properly nested.

These structures are not common, but have been experimentally verified in numerous RNAs [31]. In addition to this, these kinds of structures also appear to perform useful biological functions. For example, pseudoknots have been shown to allow frame shifting during translation of proteins [19]. In laymans terms, pseudoknots can change the way RNA is read when being translated into protein. Such frame shifting is used extensively by viruses, particularly HIV [19]. Unfortunately, the problem of finding optimal structures with pseudoknots has been shown to be NP-Complete [18].

1.2.2 Pseudoknots

In my investigation I did not consider the prediction of pseudoknotted RNA as they are uncommon. However I deem it important to understand how such structures could be integrated into the algorithms presented in this paper. As such, I have included a brisk overview of current pseudoknot prediction techniques.

Despite the problem being NP-Complete, in 1999 Rivas & Eddy [26] introduced an ingenious dynamic programming algorithm based on a new thermodynamic model encompassing pseudoknots. Their algorithm could predict a large set of pseudoknot classes using $O(N^6)$ time and $O(N^4)$ memory. They generalised the Zuker method by using a gap matrix to represent regions being considered for bonding, rather than the single continuous region used in the Zuker method. Because of its extreme space and time requirements this algorithm is used only sparingly in practice; other thermodynamic based methods for pseudoknot prediction have been formulated based on similar principals. Deogun et al. [8] described an algorithm which could handle a restricted class of pseudoknots (only those containing non-recursive pseudoknots) in $O(N^4)$ time and using $O(N^3)$ space. Shortly after which Reeder & Giegerich [23] presented an algorithm which could predict only simple recursive pseudoknots which met their ‘canonization’ criteria, and which required $O(N^4)$ time and $O(N^2)$ space. While seemingly restrictive, this did, in fact, predict a large array of pseudoknots accurately.

In recent years different approaches have been explored. In 2010, DotKnot [28] improved upon previous algorithms by using probability dot-plot guided heuristics, and an updated energy model, to predict pseudoknots. DotKnot was able to predict pseudoknots more accurately than existing methods with more frugal space and time requirements. This technique was later refined so that it could predict H-type pseudoknots and intramolecular kissing hairpins [29].

1.3 Accuracy

It is important to test and compare the accuracy of various prediction methods. As such, well established nomenclature and techniques have been developed over the history of RNA structure prediction. These methods are simple but effective. Usually accuracy is determined by comparing predicted structures to known structures. True Positives (TP) is defined as the number of base pairs which appear in both the predicted structure and the actual structure. False Positives (FP) is the number of predicted base pairs not in the true structure [17]. Similarly, False Negatives (FN) is defined as the number of base pairings in the reference structure but not present in the predicted structure [17]. Sensitivity is also called the True Positive Rate (TPR), and can be defined using the previously introduced values. I have given a mathematical definition of the TPR in Equation 1.4.

$$\frac{TP}{TP + FN} \tag{1.4}$$

Precision, also known as Positive Predictive Value (PPV), can also be calculated using these values (see Equation 1.5).

$$\frac{TP}{TP + FP} \tag{1.5}$$

RNAfold is one of the leading RNA folding algorithms, and is made available as part of the Vienna RNA package [17]. At its heart, it is an implementation of the original dynamic programming algorithm first discovered by Zuker, albeit with a more refined energy model. It is an extremely efficient implementation of this algorithm, and is also one of the most accurate in terms of sensitivity and PPV as compared to other implementations of the same algorithm [17]. When Reeder & Giegerich [23] first described their algorithm (implemented in the package pknotsRG) for pseudoknot prediction they compared it to RNAfold, and the algorithm of Rivas & Eddy (implemented in the same package and hereafter referred to as pknotsRE) [26].

Their algorithm generally had higher sensitivity than both other methods, but it is worth noting that pknotsRE was extremely close despite being based on an outdated energy model. This is possibly explained by the fact that it is a more general, and thus a more powerful algorithm. RNAfold lagged behind pknotsRE and pknotsRG in sensitivity, but executed orders of magnitude faster. Indeed, it has been shown to have excellent accuracy for smaller RNAs containing no pseudoknots while also exhibiting unrivalled computation speed [17].

1.4 Alternative Techniques

1.4.1 Soft Computing

The use of soft computing techniques has also yielded some success in the prediction of RNA secondary structure. Koessler et al. [15] modelled RNA structures as a tree of internal structures, then used artificial neural networks to recognize which of these trees appeared most RNA like. These trees are generated by constructing basic secondary structures and combinatorially merging them together to form many trees, each of which is represented as a vector of simpler trees. This vector is used as the input to the neural network.

This kind of combinatorial blending of RNA stems is also a techniques shared by genetic algorithms. Indeed, this is precisely the starting point of Van Batenburg, Gultyaev, and Pleij [34], who used a simple genetic algorithm to predict secondary structure. Their algorithm started by computing an array of all possible stems; each genome was represented as a binary string where 1 indicates a stem is in the candidate structure, and 0 indicates that it was not. Their genetic algorithm proceeded by seeding the genomes with random bits, then in a series of generation steps performed typical binary mutation, crossover and breeding, conserving and selectively breeding the fittest solutions. Fitness was defined in their algorithm as the summed total length (number of bonds) of all stems. In an improved version the summed stacking free energy reduction of all stems was used.

Unfortunately they discovered a problem with this approach: the population contained a relatively large portion of zero fitness individuals. This was because many combinations of stems are incompatible with each other, yielding impossible structures. Instead of giving these structures zero fitness, they altered their algorithm slightly to disallow crossover for stems that created an invalid structure. In addition to this, they also explored an important advantage of genetic algorithms for RNA secondary structure prediction: that of kinetic folding. Kinetic folding is the hypothesis that some RNAs, particularly large ones, have a rugged energy landscape, and because of the incremental process of transcription and folding (which happen simultaneously) become stuck in suboptimal areas during folding [33, 34].

The algorithm of Van Batenburg, Gultyaev, and Pleij simulated this process by limiting the size of stems that could contribute to a genome, and increasing this size over time until the length of the RNA was reached. This single modification to their algorithm yielded the greatest improvement in predictive power. It should also be noted that it could also predict pseudoknots, as the algorithm did not

force stems to be nested. Despite this, their approach was still less accurate than the dynamic programming approaches they compared it to. This was likely because their energy model was puerile in comparison rather than because the algorithm was flawed.

Indeed, Wiese, Deschenes, and Hendriks [36] introduced an improved genetic algorithm based on the same principles as that of Van Batenburg, Gultyaev, and Pleij which instead used an advanced energy model for fitness. They then demonstrated that it outperformed the popular dynamic programming algorithm Mfold [38] which uses a similarly complex model. Wiese, Deschenes, and Hendriks noted that as the length of RNA molecules increased the correlation between lower free energy and accuracy decreased. They concluded that this was due to the incompleteness of our current thermodynamic model of RNA folding. The algorithm of Van Batenburg, Gultyaev, and Pleij avoided this problem by simulating kinetic folding. Such a heuristic was not present in the algorithm of Wiese, Deschenes, and Hendriks.

1.4.2 Context Free Grammars

RNA sequences and their secondary structures can be represented as Context Free Grammars (CFGs). Various production rules output different internal structures (such as hairpin loops, or internal bulges), with the terminals producing the bases A, U, G, and C. This is a fundamentally different approach to those discussed previously, however it can, in fact, use the same thermodynamic energy model. Stochastic Context Free Grammars (SCFGs) can be used to encode the plausibility of structures, and thus find the most plausible looking structure using a thermodynamic model [27]. In addition, these kinds of algorithms can be trained to incorporate statistical information such as phylogenetic similarity, or machine learned parameters [27].

These kinds of methods have trouble with pseudoknots as non-nested structures are not compatible with CFGs. A notable workaround was applied by Kato, Seki, and Kasami [14], who used multiple context free grammars to model pseudoknots. However, their approach increases the time and space requirements prodigiously.

The greatest strength of context free grammar based approaches is that they can diverge from the use of free energy minimisation entirely. This is advantageous as using a physics based model—such as free energy minimization—requires a large volume of experimentally verified parameters. For this reason, many parameters are often not included in such models because they cannot be quantified empirically. The energy value of multi-branch loops, for example, is not known

and is usually guessed in modern RNA prediction algorithms. Likewise, the inter-structural interactions of hairpin loops, bulges, multi-branch loops, and internal loops has not been quantified experimentally and is thus not used as a free energy parameter. CONTRAfold [10] was one of the first SCFG based algorithms to achieve comparable performance to Zuker-like free energy minimization methods. It does away with the notion of free energy minimization altogether, and instead uses a set of trained parameters based on conditional log-linear models. CONTRAfold achieved an average prediction sensitivity higher than RNAfold, and also higher than that of Mfold [38], an RNA prediction package similar to RNAfold. It is important for the reader to remember that none of these approaches are able to model pseudoknots.

1.5 Locally Optimal Structure Prediction

DNA sequences, unlike typical RNA sequences, are very large indeed. Usually on the order of hundreds of megabytes of data. Such sequences contain DNA subsequences that code for RNAs. Functionally important RNAs typically have a recognizable secondary structure. When searching a large genome for functional RNAs, one can use a sliding window of fixed size to find locally optimal structures. This can be done by running a typical cubic time algorithm like RNAfold or Mfold at every window location. Let L be defined as the chosen window size, and N represent the length of the genome. This leads to a total complexity of $O(NL^3)$. While not prohibitive, this becomes intractable for many genomes which are typically extremely large: millions or billions of bases.

In 2004, Hofacker, Priwitzer, and Stadler [12] provided an excellent insight, and were able to lower this bound to $O(NL^2)$, making it possible to scan large genomes for interesting RNA secondary structure motifs. This was achieved by using the dynamic programming table from the previous step to quickly fill the table for the next window in quadratic time; because consecutive windows overlap, preceding information can be meaningfully used in each forward computational step. As a result it requires only a single table of size $O(L^2)$, and as such its memory complexity is only $O(N + L^2)$.

Later, in 2009, Horesh et al. [13] managed to lower the expected time bound to $O(NL)$ under the assumption that one is folding RNAs that are typical of naturally occurring sequences. This time complexity bound was experimentally verified, and their algorithm was shown to outperform that of Hofacker, Priwitzer, and Stadler.

Clearly good algorithms are available for the folding of consecutive RNA

windows. For even modest sized RNAs such algorithms are orders of magnitude faster than holistic secondary structure prediction algorithms. With the major caveat of not actually predicting a complete secondary structure, but only a set of locally optimal structures.

1.6 State of the Art: Global Optimization

Current state of the art algorithms fold RNA in a way that globally maximises their score according to some model. The Zuker algorithm, for example, finds the global minimum free energy configuration. This bias is largely due to the ‘thermodynamic hypothesis’. Anfinsen [3] presented this hypothesis as the underlying principle behind the formation of biologically active proteins. He held that protein fold into the minimum Gibbs free energy conformation in their typical biological environment; defined as the molecules physiological state: pH, temperature, ion concentration. Furthermore, through natural selection, molecules that are most likely to fold into the correct shape have evolved, and the atomic interactions of such molecules fully determine their final state. This insight has been invaluable for folding proteins and RNAs *in silico*. Despite this, it has recently become clear that methods for the prediction of RNA secondary structures have hit an upper limit in accuracy.

In her discussion of modern RNA prediction, Rivas [25] unified pseudoknot-free RNA folding algorithms. Her core observation is that all such prediction algorithms contain the same four key components: an architecture, or the production rules of a grammar; a scoring scheme, or how scores are assigned to these production rules; and the parametrization of the scoring scheme, or the specific values assigned to it. These features are referred to by Rivas, and by myself in the following discussion, as the ‘model’. The fourth and final feature is the folding algorithm used to find the best structure given the model. Here Rivas notes that the two dominant folding algorithms are interchangeable. The Cocke-Younger-Kasami (CYK) algorithm used to parse SCFGs, and the MFE algorithm based on the work of Zuker, are isomorphic for the purpose of parsing RNA grammars. Rivas additionally notes that all scoring schemes and parametrizations appear to hit an accuracy upper limit, and that complex, machine learned models are only slightly more accurate than thermodynamic models. In fact, relatively basic grammars with hundreds of parameters seem to perform almost equivalently to those with tens of thousands. While Rivas manages to unify many aspects of RNA prediction she did not recognise that all such algorithms are based on those same assumptions underpinning the thermodynamic hypothesis: they seek to globally maximise a scoring function for the final RNA secondary structure.

I propose that in large RNA molecules, local interactions are stronger than global interactions. As a result RNA molecules will misfold into a global structure made up of locally optimal structures. This is my core hypothesis. If this assumption is correct it follows that there exists a set of ‘windows’ that when folded using any reasonable model will be more accurate than the corresponding global optimum using the same model. There is already some evidence for this hypothesis. Dawson et al. [7] used variable Kuhn lengths to accurately predict RNAs containing less than 100 nucleotides. Kuhn length is the size of a segment in a polymer chain. It is a simplifying assumption that allows one to treat the entire chain as a sequence of Kuhn segments. In addition, Dawson et al. showed that the energy landscape of these RNAs, when the correct Kuhn length was applied, was funnel shaped. With such simplifying assumptions the energy landscape of RNAs is notoriously rugged with in vivo secondary structures often becoming ‘trapped’ in suboptimal states; this is most apparent in large RNAs [9].

If my hypothesis is supported, I aim to leverage the insights provided to improve the accuracy and speed of large RNA secondary structure prediction.

CHAPTER 2

Methology

2.1 Materials

2.1.1 Environment

All algorithms were implemented and tested using Ubuntu 13.10 running on an Intel i5-3210m processor with four gigabytes of RAM. The GNU C compiler version X was used to compile all C and C++ code.

2.1.2 Software

The Vienna RNA package was used as a base for all algorithms presented in this paper. This package contains many useful programs for working with RNA. The RNAfold and RNALfold functions were used in this investigation.

RNAfold is a modern implementation of Zukers folding algorithm. It predicts the MFE secondary structure of an RNA, given the primary sequence. It can also calculate the Boltzmann partition function and produce a matrix of base pair probabilities. Both of these features were used.

RNALfold implements a sliding window RNA folding algorithm that is designed to find all locally optimal secondary structures using a window of fixed size L . It is an implementation of Hofacker, Priwitzer, and Stadler's algorithm [12], and thus uses $O(NL^2)$ time and $O(N + L^2)$ space. I note here that this is not the most optimal algorithm available. Horesh et al. [13] presented an algorithm that also folds consecutive windows using a similar MFE model. However, they achieved a typical time complexity of $O(NL)$. This algorithm was not used because the implementation given by the authors is based on an older energy model taken from the Mfold package. RNALfold is based on the same energy model as RNAfold, which has been recently updated [17]. For the sake of accurate comparison to RNAfold, and to ensure a state-of-the-art energy model,

RNALfold was used.

Version 2.1.6 of the ViennaRNA package was used. The package was built from the C source code after minor modifications were made to Lfold.c and Lfold.h. The Vienna RNA makefile was used to compile the package, and the resulting RNALib static library was linked at compile time.

The Lfold.c and Lfold.h files were modified so that, when the algorithm was executed, RNALfold returned a linked list of local secondary structures and their free energy. Before modification it would instead print them to the standard output stream. The modified versions of these files are available here.

Statistical tests were performed using SOFA Statistics version 1.4.2. Regression analyses and plots were generated using Gretl version 1.9.14.

2.1.3 Testing Set

The RNA secondary structures used to test the algorithms presented in this paper were taken from the RNA STRAND database [2]. The RNA STRAND database is a free-to-use, curated collection of RNA secondary structures taken from various publicly available databases and publications. A subset of RNA structural data was extracted from the database. This subset contained only RNA structures that were marked having been verified using X-Ray Crystallography, or Nuclear Magnetic Resonance imaging. It also comprises only whole RNAs; none of the RNAs used were fragments or subsequences of larger RNA molecules. Finally, no duplicates were allowed in the test set.

2.2 Merging RNA Intervals

Running an algorithm to generate consecutive windows of size L over a RNA primary sequence of size N will produce $N - L + 1$ windows of size L . Each of these windows may contain more than one discrete structural element. For example, a window may contain two disjoint stems. After generating these windows, a post-processing step was executed in which any such disjoint secondary structural elements were split up into their own ‘intervals’. While RNALfold does this without modification, it might be necessary to encode this manually using a different sliding window algorithm implementation.

Generating a set of locally optimal secondary structure intervals is trivial once the Vienna RNA package is set up correctly. However, many of the predicted intervals overlap. To produce a valid secondary structure from these substructures,

tures, I devised a procedure to select a set of non-intersecting windows. Such a set should form the most accurate possible secondary structure when merged. Several procedures were tested.

Let W be the set of all windows to choose from, and let the set S represent the set of selected secondary structures. Also, note that every window has an associated score, which relates inversely to its free energy contribution. This is to say that a higher score is equivalent to a lower free energy.

2.2.1 Top-Down Selection

In this method, W is sorted by window size in descending order. The algorithm then examines each element of W in order. If an element of W is compatible, it is added to S . A window is compatible if it does not overlap any window in S , and if it neither contains nor is contained by any window in S .

Algorithm 1 Top-Down Selection

```

1: Let  $W$  be the set of RNA intervals
2: Let  $S = \{\}$ 
3: Sort( $W$ )
4: for all  $W$  as  $e$  do
5:   if COMPATIBLE( $e, S$ ) then
6:      $S = S \cup e$ 
7:   end if
8: end for
9: return  $S$ 

```

The algorithm (see Algorithm 1) examines every element of W exactly once, and for each element checks if it is compatible with every element in S . The naive way to check for compatibility would be to examine every element in S . This would give the algorithm a final worst case time complexity of $O(N^2)$, where N is the number of intervals, as all the intervals in W may be compatible. This is a result of the algorithm first checking S of size 0, then 1, up to $N - 1$. The closed form of this is $(N - 1)((N - 1) + 1)/2$. Additionally it requires $O(N)$ space to store S .

It is possible to lower the time complexity to $O(N \log N)$ by storing the elements of S in an Augmented Interval Tree. Querying this tree finds any elements of S that intersect the interval in question. This query takes $O(\log N)$ time and is executed N times—once for every element of W .

This is optimal, since the **Sort** operation requires $O(N \log N)$ time. Hence,

the final worst case time complexity is $O(N \log N)$ and the final worst case space complexity is $O(N)$.

2.2.2 Bottom-Up Selection

This algorithm is identical to Top-Down Selection, except that W is sorted by interval size in ascending order.

2.2.3 MFE Selection

This algorithm is identical to Top-Down Selection, except that W is sorted by score in descending order.

2.2.4 Weighted Activity Selection

The weighted activity selection problem is a generalization of the activity selection problem. Given a set of intervals, each of which is assigned a weight, the problem is to find a subset such that none of these intervals are touching one-another, and the sum of their weights is maximised. This maps naturally to the problem domain of this investigation. I used this algorithm to find the subset of compatible windows which have the minimum sum free energy, or put another way, the maximum score.

Algorithm 2 Weighted Activity Selection

```

1: Sort  $W$  by right end points
2: Let  $N = |W|$ 
3: Let  $q = \text{array}[1..N]$ 
4: for  $i = 1 \rightarrow N$  do
5:    $q[i] = \text{highest index} < i$  which is compatible with  $i$ 
6: end for
7: Let  $dp = \text{array}[0..N]$ 
8:  $dp[0] = 0$ 
9: for  $i = 1 \rightarrow N$  do
10:   $dp[i] = \max(\text{weight}[i] + dp[q[i]], dp[i - 1])$ 
11: end for
12: return  $dp[N]$ 

```

The initial sort (see line 1 in Algorithm 2) requires $O(N \log N)$ time. In the following loop (line 4) the algorithm finds the right most element in W , such that

it is compatible with i (see Figure X). This can be done using a binary search. Hence the total worst case run time of the loop is also $O(N \log N)$. The final dynamic programming array fill (line 9) requires only $O(N)$ time. As a result the final time complexity of this algorithm is $O(N \log N)$. Additionally, $O(N)$ space is required for the arrays q and dp .

2.2.5 Comparison of RNA Interval Selection Algorithms

The reader should note that all presented selection algorithms have the same $O(N \log N)$ worst case time complexity, and the same auxiliary space complexity of $O(N)$. None of these algorithms have large constant factors. Therefore, my final choice of algorithm was based solely on the accuracy of the computed secondary structure.

In order to test and compare these selection algorithms, all the windows of size 5 to 500 were precomputed for every RNA in the testing set. This is to say that the sliding window algorithm was run for a window size of 5, then 6, up to and including size 500 for each RNA. The minimum size of 5 was chosen because meaningful stems rarely form with a size fewer than five nucleotides. The maximum of 500 was due to space and time constraints. Computing these windows took several days, and used 2.5 gigabytes of RAM to store.

Each selection algorithm was run using these data. For every RNA, the single most accurate window size was chosen. Accuracy was judged to be the F1-Score, which is the harmonic mean of sensitivity and precision (positive predictive value).

2.3 Prediction With Windows

In this approach I computed sliding windows using a ‘splat’ of sizes. WAS was then done on the resulting set of computed windows.

I explored several methods for finding good splats. I present here to most successful technique, which I call ‘AB-Splat’ prediction. Initially the algorithm started at a small window size, and exponentially increased the scope of the window until a threshold was exceeded. This is explained using pseudocode in Algorithm 4.

The time complexity of this algorithm is non-trivial to deduce. The runtime of `ComputeSlidingWindows` is $O(NL^2)$, which is called one for every window of size $O(\sqrt{N})$, due to the threshold function. The sizes of windows are exponentially

Algorithm 3 Splat RNA-Prediction

```
1: Let  $RNA$  be the primary RNA sequence being folded
2: Let  $Splat$  be the set of window sizes
3: Let  $S = \{\}$ 
4: for all  $Splat$  as  $s$  do
5:    $S = S \cup \text{COMPUTESLIDINGWINDOWS}(RNA, s)$ 
6: end for
7: return  $\text{WEIGHTEDACTIVITYSELECTION}(S)$ 
```

Algorithm 4 AB-Splat

```
1: Let  $RNA$  be the primary RNA sequence being folded
2: Let  $i = a$ 
3: Let  $S = \{\}$ 
4: while  $i \leq \text{TRESHOLD}(RNA)$  do
5:    $S = S \cup \text{COMPUTESLIDINGWINDOWS}(RNA, i)$ 
6:    $i = i \times b$ 
7: end while
8: return  $\text{WEIGHTEDACTIVITYSELECTION}(S)$ 
```

expanding Thus we arrive at

$$\begin{aligned} & na^2 + n(ab)^2 + n(ab^2)^2 + n(ab^3)^2 + \dots + n^2 \\ & n^2 \left(\frac{na^2}{n^2} + \frac{n(ab)^2}{n^2} + \frac{n(ab^2)^2}{n^2} + \frac{n(ab^3)^2}{n^2} + \dots + 1 \right) \\ & O(n^2) \end{aligned}$$

Every window size produces n windows of size i .

$$\begin{aligned} & na + nab + nab^2 + nab^3 + \dots + n^{\frac{3}{2}} \\ & O(n^{\frac{3}{2}}) \end{aligned}$$

Weighted Activity selection uses $O(n \log n)$ time. As we have $O(n^{\frac{3}{2}})$ windows in the worst case, this leads to a worst case time complexity of $O(n^{\frac{3}{2}} \log n^{\frac{3}{2}})$ which is dominated by $O(n^2)$.

Zuker's algorithm requires $O(n^3)$ time and $O(n^2)$ space. The computational bottleneck is usually time, however.

A brute force method was used to test a combinatorial set of values for a and a such that $10 \leq a \leq 30$ and $1.5 \leq B \leq 4.0$. In this method all integer values for

a were attempted, and values for b were generated with a step size of 0.1. The threshold function was defined as $\min\{\sqrt{RNALength} \times 9.5, 500\}$.

To test this algorithm the test set was randomly partitioned into two sets containing an equal number of RNAs. One of these sets I shall call the training set, the other I shall call the validation set. A brute force search was done (as described in ref) on the training set, and the F1-Scores for a and b pairs was recorded; ab -pairs with highest F1-Scores were deemed the best.

These scores were then validated by running the same procedure on the validation set, then attempting to correlate the scores of ab pairs between training and validation sets. An ordinary least squares linear regression was done for correlation. The resulting model indicates if high scores in the training set are related to high scores in the testing set. This allowed me to determine if results found in the training set are valid.

The best ab -pair was used to configure the splat algorithm, which was then run on the full test set, and it's F1-Scores was recorded. These F1-Scores were then compared directly to F1-Scores produced by RNAfold, which was also run on the full test set. In addition to this, the time to execute each algorithm was recorded and compared.

Bibliography

- [1] ALBERTS, B., BRAY, D., HOPKIN, K., JOHNSON, A., LEWIS, J., RAFF, M., ROBERTS, K., AND WALTER, P. *Essential Cell Biology*, third ed. Garland Science: New York, 2009.
- [2] ANDRONESCU, M., BEREG, V., HOOS, H. H., AND CONDON, A. Rna strand: the rna secondary structure and statistical analysis database. *BMC bioinformatics* 9, 1 (2008), 340.
- [3] ANFINSEN, C. Principles that govern the protein folding chains. *Science* 181 (1973), 233–230.
- [4] BEATON, M. J., AND CAVALIER-SMITH, T. Eukaryotic non-coding dna is functional: evidence from the differential scaling of cryptomonad genomes. *Proceedings of the Royal Society of London. Series B: Biological Sciences* 266, 1433 (1999), 2053–2059.
- [5] CONDON, A. Problems on rna secondary structure prediction and design. In *Automata, Languages and Programming*. Springer, 2003, pp. 22–32.
- [6] CONN, G. L., AND DRAPER, D. E. Rna structure. *Current opinion in structural biology* 8, 3 (1998), 278–285.
- [7] DAWSON, W., TAKAI, T., ITO, N., SHIMIZU, K., AND KAWAI, G. A new entropy model for rna: part iii, is the folding free energy landscape of rna funnel shaped? *Journal of Nucleic Acids Investigation* 4, 1 (2013).
- [8] DEOGUN, J. S., DONIS, R., KOMINA, O., AND MA, F. Rna secondary structure prediction with simple pseudoknots. In *Proceedings of the second conference on Asia-Pacific bioinformatics-Volume 29* (2004), Australian Computer Society, Inc., pp. 239–246.
- [9] DITZLER, M. A., RUEDA, D., MO, J., HÅKANSSON, K., AND WALTER, N. G. A rugged free energy landscape separates multiple functional rna folds throughout denaturation. *Nucleic acids research* 36, 22 (2008), 7088–7099.
- [10] DO, C. B., WOODS, D. A., AND BATZOGLOU, S. Contrafold: Rna secondary structure prediction without physics-based models. *Bioinformatics* 22, 14 (2006), e90–e98.

- [11] HOFACKER, I. L. Rna consensus structure prediction with rnaalifold. In *Comparative Genomics*. Springer, 2008, pp. 527–543.
- [12] HOFACKER, I. L., PRIWITZER, B., AND STADLER, P. F. Prediction of locally stable rna secondary structures for genome-wide surveys. *Bioinformatics* 20, 2 (2004), 186–190.
- [13] HORESH, Y., WEXLER, Y., LEBENTHAL, I., ZIV-UKELSON, M., AND UNGER, R. Rnaslider: a faster engine for consecutive windows folding and its application to the analysis of genomic folding asymmetry. *BMC bioinformatics* 10, 1 (2009), 76.
- [14] KATO, Y., SEKI, H., AND KASAMI, T. Stochastic multiple context-free grammar for rna pseudoknot modeling. In *Proceedings of the Eighth International Workshop on Tree Adjoining Grammar and Related Formalisms* (2006), Association for Computational Linguistics, pp. 57–64.
- [15] KOESSLER, D. R., KNISLEY, D. J., KNISLEY, J., AND HAYNES, T. A predictive model for secondary rna structure using graph theory and a neural network. *BMC bioinformatics* 11, Suppl 6 (2010), S21.
- [16] LEUNG, Y. Y., RYVKIN, P., UNGAR, L. H., GREGORY, B. D., AND WANG, L.-S. Coral: predicting non-coding rnas from small rna-sequencing data. *Nucleic acids research* 41, 14 (2013), e137–e137.
- [17] LORENZ, R., BERNHART, S. H., ZU SIEDERDISSEN, C. H., TAHER, H., FLAMM, C., STADLER, P. F., HOFACKER, I. L., ET AL. Viennarna package 2.0. *Algorithms for Molecular Biology* 6, 1 (2011), 26.
- [18] LYNDSØ, R. B., AND PEDERSEN, C. N. Rna pseudoknot prediction in energy-based models. *Journal of computational biology* 7, 3-4 (2000), 409–427.
- [19] NAM, O., MORAN, S. J., STUART, D. I., GILBERT, R. J., AND BRIERLEY, I. A mechanical explanation of rna pseudoknot function in programmed ribosomal frameshifting. *Nature* 441, 7090 (2006), 244–247.
- [20] NEIDLE, S. *Principles of nucleic acid structure*. Academic Press, 2010.
- [21] NUSSINOV, R., AND JACOBSON, A. B. Fast algorithm for predicting the secondary structure of single-stranded rna. *Proceedings of the National Academy of Sciences* 77, 11 (1980), 6309–6313.

- [22] NUSSINOV, R., PIECZENIK, G., GRIGGS, J. R., AND KLEITMAN, D. J. Algorithms for loop matchings. *SIAM Journal on Applied mathematics* 35, 1 (1978), 68–82.
- [23] REEDER, J., AND GIEGERICH, R. Design, implementation and evaluation of a practical pseudoknot folding algorithm based on thermodynamics. *BMC bioinformatics* 5, 1 (2004), 104.
- [24] REUTER, J. S., AND MATHEWS, D. H. Rnastructure: software for rna secondary structure prediction and analysis. *BMC bioinformatics* 11, 1 (2010), 129.
- [25] RIVAS, E. The four ingredients of single-sequence rna secondary structure prediction. a unifying perspective. *RNA biology* 10, 7 (2013), 1185.
- [26] RIVAS, E., AND EDDY, S. R. A dynamic programming algorithm for rna structure prediction including pseudoknots. *Journal of molecular biology* 285, 5 (1999), 2053–2068.
- [27] RIVAS, E., LANG, R., AND EDDY, S. R. A range of complex probabilistic models for rna secondary structure prediction that includes the nearest-neighbor model and more. *RNA* 18, 2 (2012), 193–212.
- [28] SPERSCHNEIDER, J., AND DATTA, A. Dotknot: pseudoknot prediction using the probability dot plot under a refined energy model. *Nucleic acids research* 38, 7 (2010), e103–e103.
- [29] SPERSCHNEIDER, J., DATTA, A., AND WISE, M. J. Heuristic rna pseudoknot prediction including intramolecular kissing hairpins. *RNA* 17, 1 (2011), 27–38.
- [30] STUDNICKA, G. M., RAHN, G. M., CUMMINGS, I. W., AND SALSER, W. A. Computer method for predicting the secondary structure of single-stranded rna. *Nucleic acids research* 5, 9 (1978), 3365–3388.
- [31] TAUFER, M., LICON, A., ARAIZA, R., MIRELES, D., VAN BATENBURG, F., GULTYAEV, A. P., AND LEUNG, M.-Y. Pseudobase++: an extension of pseudobase for easy searching, formatting and visualization of pseudoknots. *Nucleic acids research* 37, suppl 1 (2009), D127–D135.
- [32] TINOCO JR, I., AND BUSTAMANTE, C. How rna folds. *Journal of molecular biology* 293, 2 (1999), 271–281.
- [33] TREIBER, D. K., AND WILLIAMSON, J. R. Beyond kinetic traps in rna folding. *Current opinion in structural biology* 11, 3 (2001), 309–314.

- [34] VAN BATENBURG, F., GULTYAEV, A. P., AND PLEIJ, C. W. An ap-programmed genetic algorithm for the prediction of rna secondary structure. *Journal of Theoretical Biology* 174, 3 (1995), 269–280.
- [35] WATSON, J. D., CRICK, F. H., ET AL. Molecular structure of nucleic acids. *Nature* 171, 4356 (1953), 737–738.
- [36] WIESE, K. C., DESCHENES, A. A., AND HENDRIKS, A. G. Rnapredictan evolutionary algorithm for rna secondary structure prediction. *Computational Biology and Bioinformatics, IEEE/ACM Transactions on* 5, 1 (2008), 25–41.
- [37] XU, Z., ALMUDEVAR, A., AND MATHEWS, D. H. Statistical evaluation of improvement in rna secondary structure prediction. *Nucleic acids research* 40, 4 (2012), e26–e26.
- [38] ZUKER, M. Mfold web server for nucleic acid folding and hybridization prediction. *Nucleic acids research* 31, 13 (2003), 3406–3415.
- [39] ZUKER, M., AND STIEGLER, P. Optimal computer folding of large rna sequences using thermodynamics and auxiliary information. *Nucleic acids research* 9, 1 (1981), 133–148.