

Secondary structure prediction of long non-coding RNA: review and experimental comparison of existing approaches

L.A. Bugnon^{1*} A. Edera¹ S. Prochetto^{1,2} M. Gerard¹ J. Raad¹
E. Fenoy¹ M. Rubiolo¹ U. Chorostecki³ T. Gabaldón^{3,4,5} F. Ariel²
L. Di Persia¹ D.H. Milone¹ G. Stegmayer¹

July 4, 2022

¹Research Institute for Signals, Systems and Computational Intelligence
sinc(*i*) (CONICET-UNL), Ciudad Universitaria, Santa Fe, Argentina.

²IAL, CONICET, Ciudad Universitaria UNL, (3000) Santa Fe, Argentina.

³Barcelona Supercomputing Center (BSC-CNS), Institute of Research in Biomedicine (IRB), Spain.

⁴Catalan Institution for Research and Advanced Studies (ICREA), Barcelona, Spain.

⁵Centro de Investigacin Biomdica En Red de Enfermedades Infecciosas (CIBERINFEC), Barcelona, Spain.

Abstract

Motivation: In contrast to messenger RNAs, the function of the wide range of existing long non-coding RNAs (lncRNAs) largely depends on their structure, which determines interactions with partner molecules. Thus, the determination or prediction of the secondary structure of lncRNAs is critical to uncover their function. Classical approaches for predicting RNA secondary structure have been based on dynamic programming and thermodynamic calculations. In the last 4 years, a growing number of machine learning (ML)-based models, including deep learning (DL), have achieved breakthrough performance in structure prediction of biomolecules such as proteins and have outperformed classical methods in short transcripts folding. Nevertheless, the accurate prediction for lncRNA still remains far from being effectively solved. Notably, the myriad of new proposals has not been systematically and experimentally.

Results: In this work we compare the performance of the classical methods as well as the most recently proposed approaches for secondary structure prediction of RNA sequences using a unified and consistent experimental setup. We use the publicly available structural profiles for 3,023 yeast RNA sequences, and a novel benchmark of well-characterized lncRNA structures from different species. Moreover, we propose a novel metric to assess the predictive performance of methods, exclusively based on the chemical probing data commonly used for profiling RNA structures, avoiding any potential bias incorporated by computational predictions when using dot-bracket references. Our results provide a comprehensive comparative assessment of existing methodologies, and a novel and public benchmark resource to aid in the development and comparison of future approaches.

Availability: Full source code and benchmark datasets are available at: <https://github.com/sinc-lab/lncRNA-folding>

Contact: lbugnon@sinc.unl.edu.ar

Supplementary information: Supplementary data are available at *Briefings in Bioinformatics* online.

1 Introduction

For decades, the sole function assigned to RNAs was to act as an information messenger between DNA and proteins. Originally, in the central dogma of biology, RNA was considered to play a secondary part in expressing inherited information as proteins. However, our recent ability to sequence entire genomes and transcriptomes served to uncover that the majority of the human genome is transcribed, although nearly 98% of it comprises non-coding regions [37, 75]. Growing evidence has linked non-coding RNAs (ncRNAs) to virtually every step of gene expression regulation, including epigenetics and spatial organization of genetic information in the cell nucleus, transcripts processing and stability, messengers translation, post-transcriptional protein modification and degradation [51]. Non-coding RNAs can be classified according to their size into two large classes [32]. On the one hand active small RNAs, less than 50 nt in length, include microRNAs (miRNAs), small interfering RNAs (siRNAs), heterochromatic

siRNAs (hetsiRNAs), Piwi-associated RNAs (piRNAs), small nuclear and nucleolar RNAs (snRNAs and snoRNAs, respectively), among others. On the other hand, long non-coding RNA (lncRNAs), larger than 200 nt, may exert their functions as long transcripts without being processed into small RNAs [79]. In particular, several studies shed light on the role of lncRNAs in diverse cellular processes, such as cell-cycle regulation in health and disease [21, 38], including cancer and diabetes [85, 94, 20, 57, 2], among others. In plants, lncRNAs have been associated with development and the dynamic response to the environment [5, 46]. At the molecular level, lncRNAs have been linked to virtually every step of gene expression regulation, including chromosome inactivation, genomic imprinting, chromatin dynamics, protein modifications and stability [56, 58]. Similar to proteins, RNA function is mainly related to its structure [24]. In particular, it should be noted that ncRNAs perform their functions through interaction with other molecules (DNA, RNA, proteins and lipids). In this sense, the lncRNA secondary structure is decisive to determine its interactome and the related functional output [72].

The RNA molecule is an ordered sequence composed of four nucleotides (nts) or bases: adenine (A), cytosine (C), guanine (G) and uracil (U), arranged in the 5 to 3 direction. The pairing of these four bases within a RNA molecule gives rise to its secondary structure. Canonical base pairs include the WatsonCrick base pairs (AU and GC) and wobble base pairs (GU), which provide higher energetic stability to the molecule [39, 50]. These base pairs often result in the formation of a nested structure, where several pairs are stacked and one or more unpaired bases form a loop. RNA secondary structure is a 2-D representation of this self-folding. Figure 1 (top) illustrates an example of the basic motifs: (i) the double stranded regions named *stems*, obtained by the stacking of two or more consecutive base-pairs; (ii) the hairpin *loop*, a single-stranded region at the end of a stem; (iii) a *bulge*, a single stranded region which interrupts a stem on one side; (iv) an *internal loop* stops a stem on either side; and (v) a single stranded region where several stems meet called a *multi-branched loop* [1]. The RNA secondary structure is usually represented in dot-bracket notation, with matching parentheses for paired bases and dots for unpaired bases (Figure 1, bottom). Additional structures called pseudo-knots can be formed when unpaired bases match with distant ones, commonly annotated with other symbols such as angle- and curly-brackets.

Even though there is currently a wide variety of publicly available ncRNA sequences, and their numbers keep growing at an ever-increasing rate [73], most of their structures remain unknown. Therefore, the efficient determination of their secondary structure is of high interest, which can be carried out by physico-chemical methods. For example, from atomic coordinates obtained from X-ray crystallography or nuclear magnetic resonance (NMR) [71, 23, 33]. However, such methods have low throughput and it is very challenging to apply them on lncRNAs not only because of the high experimental costs and resolution limits [59], but also due to their length, low abundance in in-vivo systems and the large diversity of stable structures that they can adopt. To date, the majority of structural elements found in lncRNA sequences has been determined as patterns of base pairings using a combination of chemical or enzymatic probing [89] such as PARS [34], nextPARS [64], SHAPE-seq [43] or DMS-seq [17]. These probing data can be used to aid genome-wide RNA secondary structure prediction [9]. Yet, the structures of only a tiny fraction of RNAs have been experimentally determined, limiting the understanding of this key feature upon functional outputs.

In the last decades, many methods have been developed for the computational prediction of RNA secondary structure. Furthermore very recently, the development of RNA tertiary structure prediction methods have been proposed as well [87, 7, 41]. The first proposals, dating 15 years ago [49], were based on dynamic programming and thermodynamics calculations [65, 78], identifying a structure with minimum free energy (MFE) according to the principle that RNA molecules exist in energetically stable states, like proteins [4]. Until the irruption of machine learning (ML)-based methods in the field approximately 4 years ago, prediction accuracy has remained almost unaltered. Recent works have shown that newer methods based on Deep Learning (DL) can outperform existing mainstream methods on small datasets, in terms of accuracy and applicability [92].

DL techniques first emerged as an alternative approach to structure prediction problems in proteins with AlphaFold [30] and in RNA secondary structure prediction as well [90, 67]. Compared to classical approaches, DL methods make much weaker assumptions about the thermodynamic mechanic driving RNA folding, which is based on labor-intensive experimental melting data [78]. Thus, they are more suited to detect more complex foldings [22], such as non-canonical base pairing or previously unrecognized base pairing constraints. There are many different DL proposals, which differ in their architectural design, model input-output, training data and optimization algorithms used to adjust their parameters. In general, methods treat the input RNA as a sequence of characters defined by the bases, which can be processed by different well-known computational models for text processing, such as Long Short-Term Memory (LSTM) [25] or Transformer encoders [82] that are well suited for capturing long-range interactions between nucleotides. Other approaches integrate DL techniques with thermodynamic methods to

2 RNA secondary structure prediction methods

RNA structure prediction has been classically formulated as an optimization problem, where a score is defined for every possible folding of the given RNA sequence, and the predicted folding is the one that maximizes it. The most popular approach was based on thermodynamic models [65], such as Turner’s nearest neighbor model [78]. The free energy of each nearest neighbor is calculated by summing up its free energy parameters. The free energy of an entire RNA secondary structure will be the sum of the free energy of each nearest-neighbor loop. An optimal secondary structure with MFE is calculated using dynamic programming, such as the Zuker algorithm [95].

The scores of each local element were obtained from wet-lab experiments reflecting the thermodynamics free energy theory [65]. However, the increasing availability of known RNA structures for training made it possible to successfully drift towards hybrid approaches, which are based on dynamic programming but make a fine-tuned parameter estimation based on ML [26]. Nowadays, methods for RNA structure prediction are shifting towards full ML and DL-based approaches, due to the fact that there is now larger data available for training (mostly of short sequences). A comprehensive summary of the methods that have appeared in the last 15 years in literature is presented in Table 1 and described in detail in Supplementary Material 1. The methods included in the experimental results of this study are only those available, at the time of writing this manuscript, as open access tools (already trained prediction models) in a repository or as public web servers. There are several other methods for RNA secondary structure prediction that are valuable approaches but, unfortunately, could not be included in this review due to the lack of source code or repository availability [9, 54, 83, 45, 90, 86, 53, 35, 44, 36]. These works could be included in further evaluations once the source code or web server were made available.

3 Data and performance measures

3.1 Data

Saccharomyces cerevisiae (**sce**). A large benchmark dataset was obtained from [34] where PARS was performed to characterize the secondary structure of the messenger RNAs (mRNAs) and ncRNAs of the budding yeast *Saccharomyces cerevisiae*. A total of 3,199 transcripts ranging from 71 to 8,145 bases in length were profiled. Since this dataset is composed of mRNAs and ncRNAs, and the main goal of this review is to evaluate methods on lncRNAs, it is important to distinguish the coding from the non-coding sequences. Thus, we split it into 3 specific sub-datasets: i) sce3k, with a large number of all types of sequences; ii) sce188, including only sequences with high capability of non-coding function; and iii) sce18, with only non-coding sequences. The sce3k subset has 3,023 unique sequences with length larger than 200 nt. The subset sce188 was obtained from the previous one by identifying non-coding transcripts. To this end, we used the coding potential calculator 2 (CPC2) [31], which is species-neutral and has high accuracy for long non-coding transcripts. The sce18 subset was defined by taking only those sequences from sce188 that were not previously classified as mRNA in [34]. Authors used the Vienna package [42] to fold transcripts, calculate the partition function of the structure ensemble and base pairing probabilities, with probing scores as constraints¹.

¹https://genie.weizmann.ac.il/pubs/PARS10/pars10_catalogs.html

Table 1: Methods for RNA secondary structure prediction included in this review, based on thermodynamic calculation (top) and based on machine learning (bottom).

Method	Year	Type	Repository	Web server
CONTRAFold [18]	2006	Statistical learning	http://contra.stanford.edu/contrafold/download.html	http://contra.stanford.edu/contrafold/serve.html
ContrFold [60]	2009	Statistical decision theory	https://github.com/satoken/contrfold-rna-package	http://rttools.cbrc.jp/contrfold/
ShapeKnots [15]	2010	Dynamic programming	https://rna.urmc.rochester.edu/RNAstructure.html	https://rna.urmc.rochester.edu/RNAstructureWeb/Servers/ProbKnot/ProblemKnot
ProbKnot [6]	2010	Assembling structures from base-pair probabilities	https://rna.urmc.rochester.edu/RNAstructure.html	https://rna.urmc.rochester.edu/RNAstructureWeb/
RNAstructure [55]	2010	Thermodynamics	https://rna.urmc.rochester.edu/RNAstructure.html	https://rna.thi.univie.ac.at/cgi-bin/RNAWebSuite/RNAfold.cgi
RNAfold [42]	2011	Dynamic programming	https://github.com/ViennaRNA/ViennaRNA	http://rtips.dna.bio.keio.ac.jp/ipknot/
IPknot [63, 62]	2011	Integer programming	https://github.com/satoken/ipknot	https://www.cs.bgu.ac.il/~negrebch/contextfold/
RNAshape [88]	2011	Structured-prediction learning	https://www.cs.bgu.ac.il/~negrebch/contextfold/ContextFold.1_00.zip	https://bibiserv.ebitech.uni-bielefeld.de/rnashapes
contextFold [88]	2014	Structured-prediction learning	https://bibiserv.ebitech.uni-bielefeld.de/rnashapes	https://bibiserv.ebitech.uni-bielefeld.de/rnashapes?id=knashapes.view.webserve
RNAshapes [28]	2014	Abstract shape analysis	https://bibiserv.ebitech.uni-bielefeld.de/pkiss	https://linearfold.org
pKiss [28]	2014	Abstract shape analysis	https://bibiserv.ebitech.uni-bielefeld.de/pkiss	https://linearfold.org
LinearFold [26]	2019	Dynamic programming + statistical learning	https://github.com/LinearFold/LinearFold	https://linearfold.org
LinearPartition [91]	2020	Dynamic programming + base pairing probabilities	https://github.com/LinearPartition/LinearPartition	https://linearfold.org/partition
SPOT-RNA [67]	2019	ResNet + BiLSTM	https://github.com/jaswinder Singh2/SPOT-RNA	https://sparks-lab.org/serve/spot-rna
E2Fold [10]	2020	Transformer	https://github.com/m44bio/e2fold	
rnashapes-inf [84]	2020	Bi-LSTM	https://github.com/dwflmott/rna-state-inf	
SPOT-RNA2 [68]	2021	Ensemble of deep learning models	https://github.com/jaswinder Singh2/SPOT-RNA2	https://sparks-lab.org/serve/spot-rna2
MXfold2 [61]	2021	Deep learning + thermodynamic parameters	https://github.com/keio-bioinformatics/mxfold2	http://www.dna.bio.keio.ac.jp/mxfold2/
2dRNA-Fold [47]	2021	Deep learning + reinforcement learning	https://github.com/Urinx/2dRNA-Fold	http://biophy.hust.edu.cn/new/2dRNA
UFold [22]	2022	U-net	https://github.com/uci-cbecl/UFold	https://ufold.ics.uci.edu

Table 2: Dataset of curated and validated lncRNAs.

Name	Species	Length	Year	Measurement technology
NORAD#1 [12]	<i>H. sapiens</i>	1,903	2021	nextPARS
NORAD#2 [12]	<i>H. sapiens</i>	1,862	2021	nextPARS
NORAD#3 [12]	<i>H. sapiens</i>	1,614	2021	nextPARS
CYRANO[29]	<i>H. sapiens</i>	4,419	2020	SHAPE
MEG3 [80]	<i>H. sapiens</i>	1,595	2019	SHAPE
RepA [40]	<i>M. musculus</i>	1,630	2017	SHAPE + chemical probing
PAN [74]	<i>Human gammaherpes virus 8</i>	1,077	2017	SHAPE-MaP
XIST [69]	<i>M. musculus</i>	17,779	2016	SHAPE-MaP
lincRNAp21sense [11]	<i>H. sapiens</i>	311	2016	SHAPE
lincRNAp21antisense [11]	<i>H. sapiens</i>	303	2016	SHAPE
HOTAIR [70]	<i>H. sapiens</i>	2,154	2015	SHAPE + chemical probing
MALAT1 [8]	<i>H. sapiens</i>	8,415	2014	SHAPE
ROX2 [27]	<i>D. melanogaster</i>	573	2013	SHAPE+PARS

Curated set of lncRNAs. As indicated in [48], an informative benchmark requires a variety of high-quality structures. We have compiled a novel dataset including curated structures of lncRNAs that were determined by experimental analyses, for which the structures have been carefully studied and validated through a variety of biochemical approaches by other authors. We have included those available from literature and those mentioned in a very recent review on lncRNAs [59]. The curated lncRNAs included in our study are shown in Table 2, where the name, species, length, year of discovery and probing methodology are indicated in the columns, for each test lncRNA in the rows. These species and datasets were included according to their full availability. Those lncRNA are publicly available and their secondary structures have been carefully studied, based on biochemical approaches and then curated manually. We have used in this study only those lncRNA from literature for which all the information needed for the comprehensive experimental comparisons of this review is available. For each lncRNA, its reference structure is represented in two different ways. In the first place, the classical dot-bracket representation, in which the nested base-pairs 2D conformation was obtained with a classical method, using the probing scores as constraints. The other reference representation consists of the probing score per base, which indicates the pairing probability of each nucleotide in the sequence. It should be noted that we used the original secondary structures as provided by their corresponding authors for each curated lncRNA used in this study. More details on how the secondary structures were obtained with probing constraints can be found in Supplementary Material 2.

Since each experimental technique provides a different distribution of score values, a normalization² step was performed to place all the scores within the $[0, 1]$ interval, where 0 is unpaired, 1 is paired, and values around 0.5 are uncertain. First, SHAPE scores with values $p < 0.0$ or with no data were labeled with -999 (and then ignored in all subsequent processing). Similarly, SHAPE scores with $p > 1.0$ were clipped to 1.0. Then, given that SHAPE scores are inverted (i.e. values close to 0 indicate pairing), the standardization was done with $\tilde{p} = 1 - p$. In the case of PARS and nextPARS, scores were normalized by using $\tilde{p} = (p + b)/(2b)$, where b is the maximum value for the score. For PARS $b = 7$ and for nextPARS $b = 1$.

3.2 Performance measures

Classical measures. In this case, the focus of performance measures is on the accuracy of predicted base pairs in comparison to a reference structure [48]. Pairs that are both in the prediction and in the reference structure are true positives (TP), while pairs predicted but not in the true structure are false positives (FP). Similarly, a pair in the reference structure that is not predicted is a false negative (FN), and a pair that is neither predicted nor in the true structure is a true negative (TN). To fully characterize the successes and failures of structure prediction, the F_1 score is defined as

$$F_1 = \frac{2 TP}{2 TP + FP + FN}. \quad (1)$$

These classical measures were calculated with the scorer program in RNAstructure package³.

²This is detailed in the source code, results section

³<http://rna.urmc.rochester.edu/Releases/>

Mean absolute similarity (MAS) score. We propose a new score that can reflect the similarity between the nextPARS/PARS/SHAPE probing data and the predicted dot-bracket structures. Given an RNA sequence, this score takes into account the binary paired-unpaired state of each nucleotide n in the predicted structure $b(n) \in \{0, 1\}$, and the corresponding normalized probing score $p(n) \in [0, 1]$. The paired, unpaired and average similarities are defined as

$$s^+(b, p) = 1 - \frac{1}{|M^+|} \sum_{n \in M^+} |b(n) - p(n)|, \quad (2)$$

$$s^-(b, p) = 1 - \frac{1}{|M^-|} \sum_{n \in M^-} |b(n) - p(n)|, \quad (3)$$

$$s(b, p) = 1 - \frac{1}{|M|} \sum_{n \in M} |b(n) - p(n)|, \quad (4)$$

where the sets of paired and unpaired nucleotides are

$$M^+ = \{n : p(n) \geq 1/2 + \epsilon/2\}, \quad (5)$$

$$M^- = \{n : p(n) < 1/2 - \epsilon/2\}, \quad (6)$$

and $M = M^+ \cup M^-$. In this definition, ϵ is the uncertainty level around the undetermined/no-measure score (1/2). That is, all the nucleotides with scores in the range $[1/2 - \epsilon/2, 1/2 + \epsilon/2]$ are ignored in the similarity measure.

As the normalized probing scores are distributed in the $[0, 1]$ range, the maximum/minimum value of the average similarities will generally not be 1/0. This is because predicted structures are binary, preventing perfect matchings with the real-valued probing scores. Therefore, a re-scaling step must be done to take into account the best and the worst binary structures,

$$b_b(n) = \begin{cases} 1 & \text{if } p(n) \geq 1/2 + \epsilon/2 \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

$$b_w(n) = 1 - b_b(n). \quad (8)$$

Using the average similarities for these extreme cases, the normalized MAS score is defined as

$$MAS(b, p) = \frac{s(b, p) - s(b_w, p)}{s(b_b, p) - s(b_w, p)} \quad (9)$$

This score was tested in several cases, which are provided as an interactive notebook in the source code repository (results section). It should be noted that the MAS evaluates the similarity between predictions and experimentally obtained probing scores, which aims to characterize each nt with a continuous certainty level of paired/unpaired. Unlike this, the F_1 score compares the hard predictions to a dot-bracket structure obtained with a specific prediction software. Since probing scores are only experimentally obtained, the MAS is less biased than F_1 towards a specific computational approach. Moreover, MAS takes full advantage of the information available in the continuous levels of the probing scores, instead of using a hard pairing representation.

In order to perform a statistical analysis of results, a Friedman test and post-hoc Nemenyi test for the Critical Difference (CD) diagram were used [16]. The Friedman test is a non-parametric alternative to the ANOVA used to determine whether there is a statistically significant difference between the means of three or more methods tested with the same datasets. The Friedman test uses the null hypothesis that the average performance for each method is the same; and the alternative hypothesis that at least one method average is different from the rest. If the p-value obtained is less than 0.05, the null hypothesis can be rejected. In this case, the Nemenyi post-hoc test can be performed to determine statistically which methods show different means. The Nemenyi post-hoc test returns the p-values for each pairwise comparison of means, and the CD diagram shows the pairs having an average rank that is significantly different.

4 Results

4.1 Large experimental comparison on RNA secondary structure prediction

Figure 2 reports the results for the prediction methods evaluated on the largest dataset (sce3k). The figure shows violin plots with the median on the y-axis (blue diamond) of F_1 (top) and MAS score (bottom)

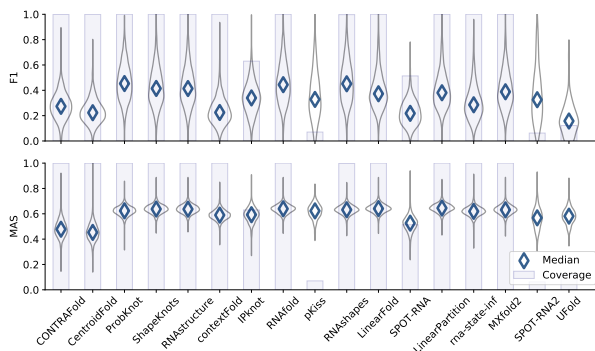


Figure 2: Comparison of performance for the RNA secondary structure prediction methods on the sce3k dataset. F_1 (top), MAS score (bottom).

for all the comparative methods in the x-axis. Methods with F_1 lower than 0.10 were not included in the analysis (E2Efold mostly predicted pseudoknots and 2dRNA-fold took a very large time to run for each sequence and could predict only a very small number of sequences). The figure also shows the corresponding coverage (gray bar in the background) for each method, i.e. the percentage of sequences from the full dataset that was effectively predicted by each method. It can be seen that most methods provided predictions for all 3,023 sequences. It should be mentioned here that we did not use the probing scores as constraints for the RNAfold predictions because this is the reference (and therefore always obtains $F_1 = 1.0$). The best methods achieved a median F_1 score of around 0.50: ProbKnot, RNashapes, RNAfold, ShapeKnots and RNAstructure. This result for such classical methods is expected considering that the reference dot-bracket for this dataset was actually obtained with RNAfold. Furthermore, the other methods with high performance share the same approach (thermodynamic modeling and dynamic programming), improving RNAfold or using it directly as part of the method. The ML based methods achieved a very close performance to the thermodynamic methods, with F_1 score around 0.40. The least performing method according to this metric achieved a F_1 score of around 0.20.

With respect to the MAS score, when looking at the similarity between experimental probing data and the predictions, most methods achieved more stable predictive performances (low variance) with median values around 0.60. According to this measure, which is not biased towards any particular computational method used to derive the reference dot-bracket, all methods (classical and ML based) have a median performance higher than 0.50. Also, it can be seen that the MAS scores distribution for each method is less dispersed than the corresponding F_1 scores distribution, showing a more consistent measurement across the different sequences in the dataset. According to MAS, ShapeKnots, RNAstructure, RNAfold, RNashape, LinearFold, LinearPartition and the ML-based MXfold2 have scores around 0.60 (and full coverage). SPOT-RNA, SPOT-RNA2 and UFold have a MAS of ≈ 0.50 .

In order to provide a statistical analysis of results, we used Friedman test and CD with post-hoc Nemenyi test ($\alpha = .05$). The Friedman test showed that the differences in the F_1 and MAS scores distributions are statistically significant ($p < 1E-15$). The CD diagram is shown in Figure 3. Here, the corresponding CD diagram for F_1 indicates that ProbKnot, RNashapes and RNAfold are the best methods for this dataset. Then there is a second group of statistically similar methods, one of which is based on ML (MXfold2). Very interestingly, the rest of the methods are well separated (different) among them. Regarding MAS, the best method is LinearPartition, well-separated and followed by classical methods and MXfold2, which are not statistically different among them.

4.2 Experimental comparison focused on lncRNAs

Figure 4 shows the comparison of prediction methods on the sce188 dataset, composed of 188 transcripts with lengths between 200 and 1,301 bases. The figure reports violin plots of the F_1 (top) and MAS scores (bottom) for all the comparative methods (x-axis). Methods with F_1 lower than 0.10 were not included. The figure also shows the corresponding coverage (gray bar) for each method. In this case, almost all methods could predict 100% of the sequences, except for UFold whose predictions are restricted to a maximum of 600 nt per sequence. The best methods according to F_1 are, again, those sharing the same classical approach as the method used in the reference structures. This result, as stated before, is expected because the reference structure was obtained with RNAfold. However, there is a larger difference in performance among methods. Many classical works have performance lower than $F_1 = 0.40$.

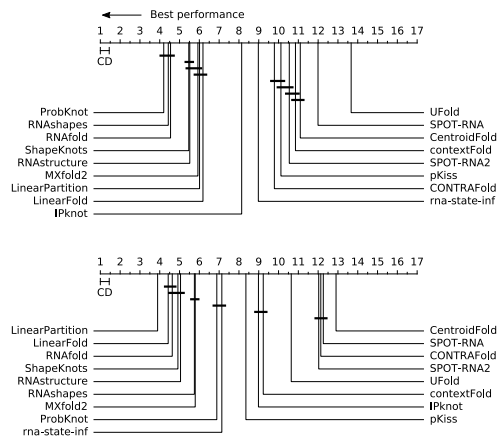


Figure 3: Critical difference diagram for the RNA prediction methods on the sce3k dataset for F_1 (top) and MAS score (bottom).

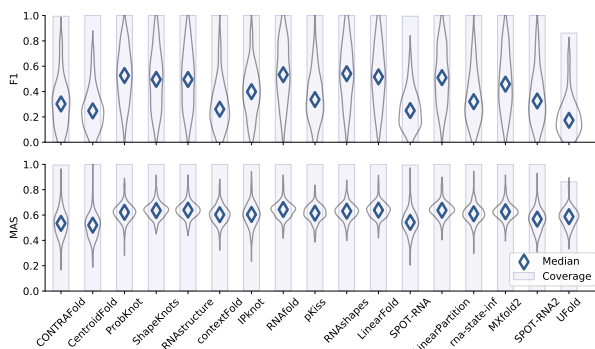


Figure 4: Comparison of performance for the RNA secondary structure prediction methods on the sce188 dataset. F_1 (top), MAS score (bottom).

Furthermore, a larger variance in each method is observed. The ML methods SPOT-RNA and MXFold2 have $F_1 = 0.40$, not that far from the best ones ($F_1 \approx 0.50$). The rest of the methods have F_1 between 0.17 and 0.40.

The MAS score shows that most methods are very close in performance, between 0.50 and 0.60. Unlike the results obtained from the sce3k dataset, the slightly wider dispersion observed here in the plots indicates that there is more heterogeneity in the predictions of sce188, which only includes transcripts predicted as lncRNAs. It can be noticed that when the MAS score is used here for measuring performance, since it is not biased towards the RNAfold-based reference, it turns out that the best classical and ML based methods are all equally good for predictions. The best classical method is RNAfold and the best ML based method is MXFold2, both with the same score around 0.60. The Friedman test indicates that differences in the F_1 and MAS score are statistically significant ($p < 1E-15$). However, given the smaller number of sequences, the statistical power of the CD method is lower than that of the sce3k dataset (see Supplementary Material 3, Figure S1).

Given the close results achieved by the methods in this dataset, the similarity of each prediction was analyzed in more detail. To this end, for each sequence the structures predicted by a pair of methods were compared in terms of which nucleotides were predicted as paired or unpaired. The average rate over the predictions obtained for each pair of methods is shown with a pair-wise heatmap in Figure 5. These comparisons show that the most similar methods are RNAstructure and ShapeKnots, which is very reasonable since the last one is based on the same thermodynamic model. This high similarity is in agreement with the previous results of F_1 and MAS score, as well as with the critical differences where both methods are similarly ranked. The second most similar methods are LinearFold and RNAfold, which is expected since the first method was proposed precisely to improve the computational performance (from cubic to linear time) of RNAfold without modifying the results. Interestingly, it can be seen that the structures predicted by CONTRAFold, CentroidFold, IPknot, SPOT-RNA and SPOT-RNA2 very

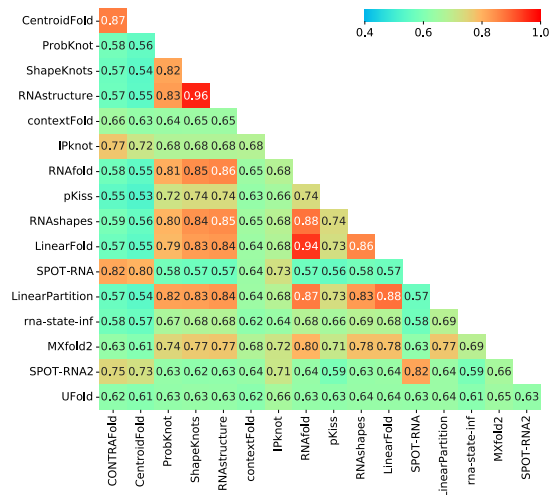


Figure 5: Paired/unpaired prediction rate between methods on the dataset sce188. Each value of the heatmap shows how similar the predictions of two methods are on average.

frequently display high similarity between them, and are different in comparison to the other methods.

Figure 6 shows the comparative results for the sce18 dataset. Surprisingly here, all methods achieved a noticeably better performance. In sce18, most methods have very high median F_1 (around 0.80), with RNAfold and LinearFold showing a median $F_1 = 0.86$, very closely followed by LinearPartition and ProbKnot ($F_1 = 0.83$), RNAstructure and ShapeKnot ($F_1 = 0.82$), MXfold2 ($F_1 = 0.80$) and SPOT-RNA2 ($F_1 = 0.77$). Importantly, these last two methods are ML-based, suggesting that they are able to achieve good predictive performance for lncRNAs. However, there is a large dispersion in the F_1 plots, hinting at a high heterogeneity in the predictions for the sce18 sequences. The other methods achieved lower scores, which are between $F_1 = 0.50$ (UFold) and $F_1 = 0.74$ (SPOT-RNA). Regarding the MAS score, all methods achieved median scores between 0.60 and 0.70, with more compact predictions and less dispersion than F_1 . Again like in the previous experiment, the best classical method is RNAfold and the best ML based method is MXFold2, both with approximately the same score, around 0.70. This indicates that all methods, no matter the approach, can reach dot-bracket structures with high correlation to the probing data. For this dataset, these results indicate that ML-based methods achieved performances very close to the classical approaches, and which are also considerably high for lncRNAs. This is an impressive result for ML-based models in comparison to very well-known and established methods for RNA structure prediction, making them very competitive for this task (Supplementary Material 3, Figure S2).

In order to determine differences in the computational cost of the methods, the running time of each method has been calculated based on the sequences in the sce18 dataset (Supplementary Material 3, Figure S3). For this calculation, all methods were run in the same hardware conditions. These results show how the sequence length clearly impacts on the computing time of each method. Almost all methods increase time for longer transcripts with the same law, showing different exponential behavior, with more cost for pKiss and ShapeKnots.

Finally, in order to understand why all methods had such high performance on the sce18 dataset, we have measured the maximum length of the hairpins present in the reference structures. The average number of hairpins per sequence was 22.32 for sce3k, 7.48 for sce188 and 5.72 for sce18, while the average of the maximum hairpin length per sequence (normalized by the sequence length) was 0.08 for sce3k, 0.16 for sce188 and 0.23 for sce18. Interestingly, sce18 sequences have the fewest but longest hairpins compared to those found in the sce188 and sce3k. We hypothesize that the sce18 secondary structures are more stable and well-formed, and thus are easier to predict by any method. A possible explanation may be that lncRNAs are more stable since their function is more dependent on the structure, compared to mRNAs. This is in line with previous findings showing that mRNAs are poorly structured [34].

4.3 Prediction analyses on a curated list of well-characterized lncRNAs

To further evaluate the methods, we have tested the RNA structure prediction with a set of lncRNAs whose structures have been experimentally validated (details in Table 2). Figure 7-(left) reports violin plots with MAS score, for each method in the x-axis. These results show that most methods provide

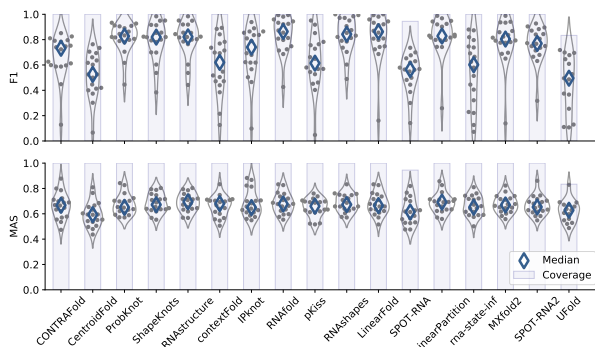


Figure 6: Comparison of RNA prediction methods on the sce18 dataset. F_1 (top), MAS score (bottom).

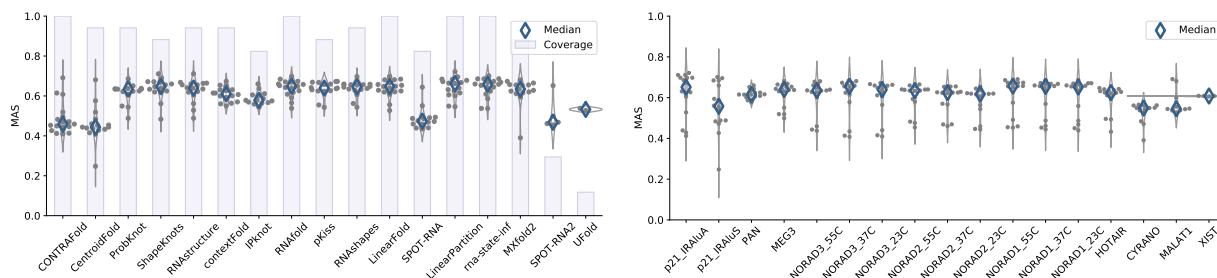


Figure 7: Comparison of MAS score performance for the RNA secondary structure prediction methods on the curated lncRNAs dataset. Results grouped by prediction method (left); each dot represents a lncRNA. Results grouped by sequence (right); each dot represents a prediction method.

predictions with high similarity to probing scores, with median values larger than 0.60, and for the least performing methods the median similarities are around 0.50 for CONTRAFold, CentroidFold, SPOT-RNA, SPOT-RNA2 and Ufold.

In order to analyze the differences in the prediction of the curated lncRNAs, Figure 7-(right) shows violin plots with the MAS score of the methods for each lncRNA in the x-axis. In this figure, the lncRNAs are ordered according to their lengths, from short to large. The median similarities are around 0.60 for all lncRNAs, no matter the length of the lncRNA to predict. Indeed, for the largest sequence (XIST) the median MAS score is 0.61 and all methods are very close to this median. This indicates that actually all methods and all approaches are capable of predicting the XIST, and also PAN, pairing patterns considerably well as measured by probing methods. It is also interesting to notice that MAS scores are very similar for the three NORAD fragments at the three different temperatures, which is expected given that the fragments resemble each other and have a similar structural function [12]. In comparison, the F_1 violin plots for this dataset (see Supplementary Material 3, Figure S4) show a large dispersion among methods predictions, for example the performance ranges from 0.10 to 0.80 for the same method depending on the lncRNA to predict. This variance makes it hard to distinguish among methods.

4.4 Detailed analysis of a curated lncRNA: NORAD

For an in-detail comparison across methods, we chose NORAD as a leading case since it is dysregulated in various types of cancer and the function is largely mediated by its structure. Figure 8 shows the results for the prediction of the secondary structure of NORAD#1 at 37 °C or all the methods reviewed here. NORAD#1 is a fragment of 1,903 bases in length within NORAD, whose reference structure has been profiled in high detail [12] using nextPARS pipeline [13]. The reference structure was built with RNAstructure and the experimentally determined nextPARS scores as constraints. For this sequence, 2dRNA-Fold could not be run because of not enough RAM, and E2EFold and Ufold because of restrictions on sequence length. The figure shows in bars the resulting predictive performances based on F_1 and MAS score. In addition, it also shows the normalized paired MAS, indicated with '+', and the normalized unpaired MAS, indicated with '-'.

The analysis of the F_1 indicates that the best methods are LinearPartition and pKiss, but with low

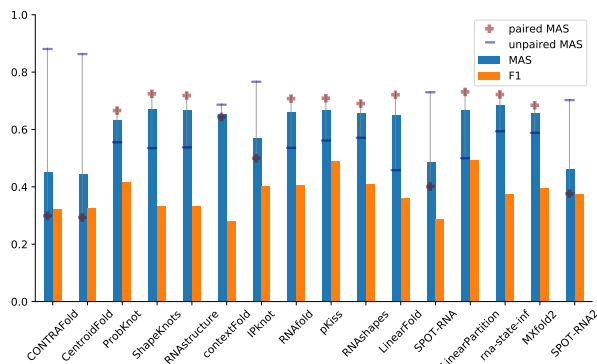


Figure 8: Comparison of F_1 and MAS scores for the RNA secondary structure prediction methods on NO-RAD#1 at 37°C.

F_1 , closely followed by ProbKnot and RNAstructure. It has to be noticed that the reference structure was built with RNAstructure + nextPARS scores as constraints. According to the MAS score, most methods were capable of correctly predicting the 70% of the pairing in the sequence. Here rna-state-inf is the best method, closely followed by ShapeKnots, RNAstructure, LinearPartition, RNAfold, pKiss and MXFold2, which are all equally good performing. The least performing methods reached a 0.45 of MAS score. The paired and unpaired MAS scores are good indicators of the capability of methods of correctly predicting both types of pairing in the sequence. For example, all winning methods here are good predictors of paired nucleotides, while methods with the lowest performance predict simpler structures and mostly identify unpaired nucleotides.

For a more comprehensive comparison among methods, we performed a deeper analysis of the RNA structures predicted by the top-3 best methods and the one with the lowest performance (Figure 9). For this analysis, the package *draw_rna*⁴ was used, which plots the structures and colors them in a range from blue to yellow (5' → 3', respectively). For visual comparisons, the same color indicates the same part of the sequence. It is very interesting to see that, according to F_1 (top) the 2nd best method is pKiss, but it can be clearly seen that the resulting structure is quite different from the reference one. Moreover, the method with the lowest F_1 (top, right) seems to be more similar to the reference one than the 2nd best method. Instead, the MAS score (bottom) provides a better ranking: all 3-top methods are visually very similar to the reference one; and the prediction of the method with the lowest performance (bottom, right) is very different to the reference.

For the top-best method in Figure 10 we used the software Circlecompare⁵ that compares two RNA structures (reference and prediction) for the same sequence by showing the nucleotides in a circle. It uses the following color scheme: green for pairs present in both predicted and reference structure, red for pairs present in predicted structure only, and black for pairs present in reference structure only. It can be seen that there are more bases in green for LinearPartition, which is reflected by its high F_1 score. It also can be seen that rna-state-inf proposes different base pairs, but most of the paired nucleotides in the reference are paired in the prediction, which explains its high MAS score. Also, both methods predict a high number of pair bases that are not part of the reference structure. Moreover, most of the long-range pairings are different from those from the reference, while local structure is more preserved.

The sequence of NORAD contains 12 recognizable and sequence-similar NORAD repeat units (NRUs) that originated by tandem duplication at the rise of mammals and still share sequence homology [77]. The NRUs are ≈300 nt in length. Most NRUs contain one or two binding sites for the two homologs of Pumilio (Pum) in mammals. NORAD-regulated Pum targets are enriched in genes involved in cell division, mitosis and chromosome number instability [76]. The individual NRUs can be studied independently, facilitating isolation of specific interaction partners. Figure 11 shows the comparison of MAS score performance for the RNA secondary structure prediction methods of NRU1, NRU2, NRU3 and NRU4 within NORAD#1 at 37°C. It can be seen that NRU1, NRU2 and NRU4 are hard to predict for all methods, reaching at best a similarity around 0.50. Besides, NRU 3 is predicted by most methods with scores near 0.80. The best one here were LinearFold, RNAfold, RNAshapes and pKiss.

Several conserved elements, including a small and a larger hairpin are peculiarly found in some NRUs and not others[76]. Thus, a possible explanation between performance of the methods for the different

⁴https://github.com/DasLab/draw_rna

⁵<https://rna.urmc.rochester.edu/Text/CircleCompare.html>

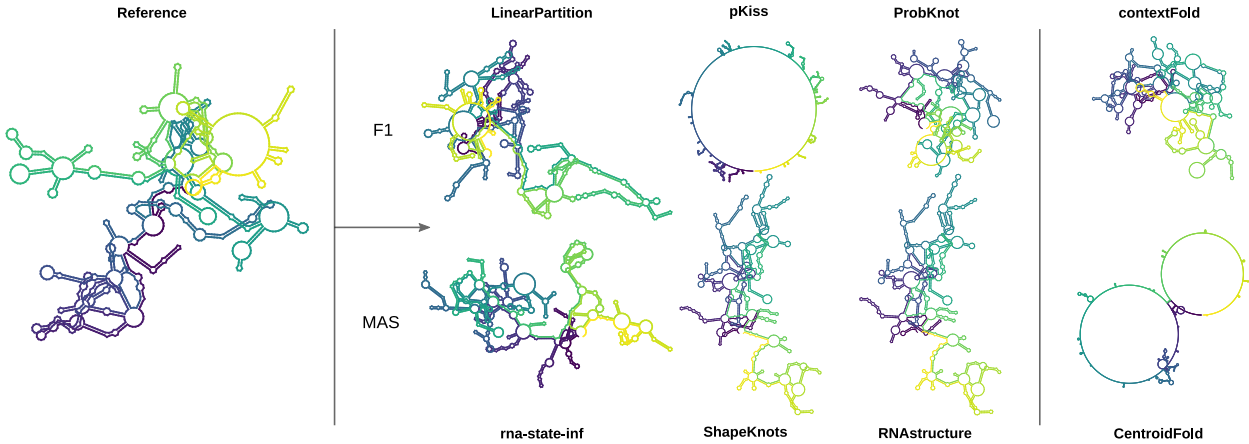


Figure 9: In-depth analysis of structure predictions for NORAD#1 at 37°C. Plots visualize the reference (left), the predictions of the top-3 methods (middle) and last prediction in the ranking (right) according to F_1 (top) and MAS (bottom). Same color indicates the same part of the sequence from 5' (blue) to 3' (yellow).

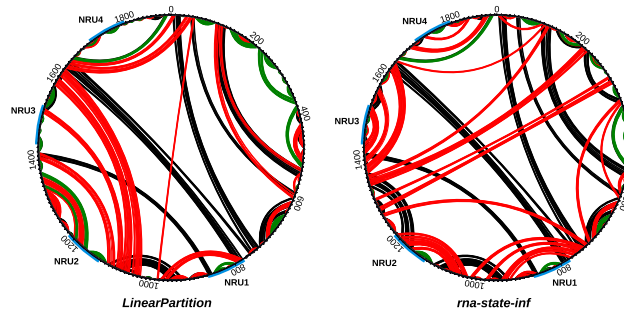


Figure 10: Circular plot for the best method in predicting NORAD #1 at 37°C according to F_1 (left) and MAS score (right). Pairs present in both predicted and reference structure (green). Pairs present in predicted structure only (red). Pairs present in reference structure only (black).

NRUs could be that, although most NRUs fold independently, there are occasional long range interactions. In the case of NRU1, it was shown that there are inter-NRU interactions between NRU1-10 [93]. Furthermore, it seems that NRU2 and NRU4 fold mostly independently, but they also have small interaction with regions outside them [12]. Instead the NRU3 folds completely independently, thus it could be locally more structured than the rest of the NRUs and that is why it is better predicted by most methods.

5 Conclusions

In this study we provided a comprehensive review and experimental comparison of classical methods as well as the most recently proposed. We used data sets validated by chemical probing methods, including a novel benchmark with a compilation of well-characterized lncRNAs from different species. Thus, our study also provides a novel public benchmark to aid in the development and comparison of future approaches.

Most available secondary structures published were obtained with computational prediction methods based on dynamic programming and thermodynamics models, which sets upis an unfair basis of comparison for methods based in any different approach. Therefore, we proposed a novel score to assess the performance of methods exclusively based on the chemical-probing data used for profiling RNA structures. The MAS evaluates the similarity between pairing predictions and probing scores, differently from the F_1 score that compares the predictions to a dot-bracket structure obtained with a specific software. Since probing scores are experimentally obtained, the MAS is less biased than F_1 towards a specific

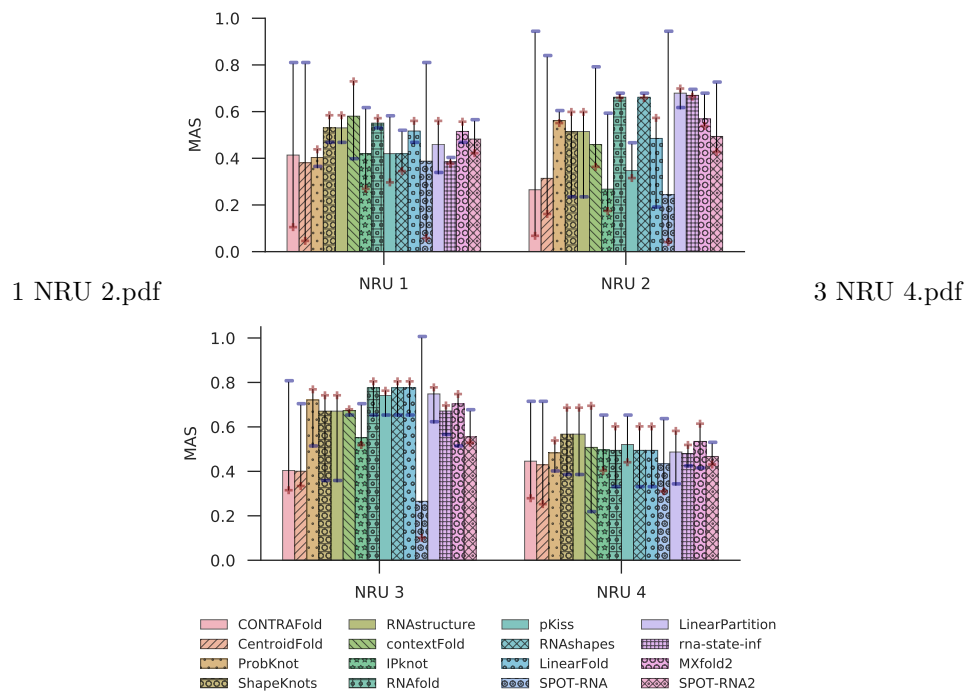


Figure 11: Comparison of MAS score performance for the RNA secondary structure prediction methods of NRU1, NRU2, NRU3 and NRU4 within NORAD #1 at 37°C.

computational approach and thus it is more faithful to the underlying data.

In summary, after the comparative results obtained in this study, it can be stated that the strength of classical methods is that they can guarantee a 50%-70% prediction performance in any type of RNA, even long non-coding transcripts. However, their weakness is that in the last 20 years, such modeling approaches have improved their computational performance but then, their results have stagnated. On the other hand, ML based efforts have been focused on reaching the performance of classical methods, not being able yet to clearly outperform them for long RNAs. Therefore, future research directions should attempt to increase this ceiling performance, probably by developing hybrid approaches between thermodynamic and ML/DL based models, since ML has the capability of boosting existing performance by modeling folding situations that have not been yet properly addressed with classical approaches.

Key Points

- We provide a comprehensive comparative assessment of existing methodologies for lncRNA secondary structure prediction in the last 15 years.
- Each method has been explained and compared experimentally on the same data set of curated lncRNAs.
- A new unbiased measure of performance is provided based only on the chemical-probing data.
- This study will help the bioinformatics community in the development and comparison of future approaches for lncRNA secondary structure prediction.

Acknowledgements

Authors would like to thank Dr. Jessica Brown, Dra. Sztuba-Solinska, Dra Anna Pyle and Rafael Araujo Tavares, Dra. Alisha Jones, Dr. Marco Marcia, Dr. Kevin Weeks and and Dr. Ilik for providing the reference data for the lncRNA experimentally validated.

Funding

This work was supported by ANPCyT (PICT 2018 3384, PICT 2018 2905, PICT 2019 3420) and UNL (CAI+D 2020 115). Researchers from sinc(i) and IAL are collaborating in the framework of the Program Science and Technology against Hunger, supported by the Argentinian Ministry of Science, to study and develop ncRNAs as exogenous bioactive molecules in agriculture. UC was funded by MICINN (IJC2019-039402-I). The work used computational resources from the Pirayu cluster, acquired with funds from the Santa Fe Science, Technology and Innovation Agency (ASACTEI), Project AC-00010-18, Res. No. 117/14. This equipment is part of the National High Performance Computing System of the Ministry of Science and Technology of Argentina. We also acknowledged the support of NVIDIA Corporation for the donation of GPUs used for this research.

Conflict of Interest: none to be declared.

Biographies

L.A. Bugnon is an Assistant Researcher in the Bioinformatics lab at sinc(i) institute, National Scientific and Technical Research Council (CONICET), and Assistant Professor at Universidad Nacional del Litoral (UNL), Santa Fe, Argentina. He works in machine learning research applied to bioinformatics.

A. Edera is a Postdoctoral Researcher in the Bioinformatics lab at sinc(i) institute. His research includes computational aspects of biological problems with focus on plants.

S. Prochetto is a Postdoctoral Researcher in the Evolution and Development lab at IAL institute, CONICET, and works closely with the Bioinformatics lab at sinc(i). His research is focused on the evolution of photosynthesis using a bioinformatics approach.

M. Gerard is an Assistant Researcher at sinc(i), CONICET, and Teaching Assistant in the Department of Informatics at UNL. His research interests include bioinformatics, machine learning and swarm intelligence.

J. Raad is a Doctoral Researcher in the Bioinformatics lab at sinc(i) institute, CONICET, and Adjunct Professor at UNL. His research interests include bioinformatics and machine learning.

E. Fenoy is a bioinformatician with experience in immunoinformatics and proteomics. He currently holds a postdoctoral position at the Research institute for signals, systems and computational intelligence, sinc(i), CONICET.

M. Rubiolo is an Assistant Researcher at sinc(i), CONICET, and Teaching Assistant at UNL and Adjunct Professor at UTN. His research interests include machine learning, data mining and bioinformatics.

U. Chorostecki is a postdoctoral researcher in the Comparative Genomics group jointly affiliated to the Biomedical Research Institute (IRB) and the Barcelona Supercomputing Centre (BSC), at Barcelona (Spain). He is interested in the study of the relationship between structure, function, and evolution in lncRNAs.

T. Gabaldón is ICREA Research Professor and Group Leader of the Comparative Genomics group at the Barcelona Supercomputing Centre (BSC) and the Institute for Research in Biomedicine (IRB) at Barcelona (Spain). He is interested in organismal and molecular evolution.

F. Ariel is Professor in the Faculty of Biochemistry at UNL, and Independent Researcher at the IAL Institute, CONICET. He is the leader of the Epigenetics and noncoding RNAs lab at IAL. His current research interest involves plant long noncoding RNAs in the regulation of genome topology dynamics.

L. Di Persia is Professor in the Department of Informatics at UNL, and Independent Researcher at the sinc(i) institute, CONICET. His research interests include biomedical signal processing, machine learning and Bioinformatics

D.H. Milone is Full Professor in the Department of Informatics at UNL and Principal Research Scientist at sinc(i), CONICET. He was founder and first Director of the sinc(i) institute. His research interests include statistical learning, signal processing and neural computing, with applications to biomedical signals and bioinformatics.

G. Stegmayer is Professor in the Department of Informatics at UNL, and Independent Researcher at the sinc(i) institute, CONICET, Argentina. She is the Leader of the Bioinformatics lab at sinc(i). Her current research interest involves machine learning, data mining and pattern recognition in bioinformatics.

References

- [1] Avinash Achar and Pål Sætrom. RNA motif discovery: a computational overview. *Biology Direct*, 10(1):1–10, October 2015.
- [2] E. Anastasiadou, L. Jacob, and F. Slack. Non-coding rna networks in cancer. *Nat Rev Cancer*, 18(1):5–18, 2018.
- [3] Mirela Andronescu, Vera Bereg, Holger H Hoos, and Anne Condon. RNA STRAND: The RNA secondary structure and statistical analysis database. *BMC Bioinformatics*, 9(1):1–10, August 2008.
- [4] Christian B. Anfinsen. Principles that govern the folding of protein chains. *Science*, 181(4096):223–230, July 1973.
- [5] Federico Ariel, Natali Romero-Barrios, Teddy Jégu, Moussa Benhamed, and Martin Crespi. Battles and hijacks: noncoding transcription in plants. *Trends in Plant Science*, 20(6):362–371, June 2015.
- [6] Stanislav Bellaousov and David H. Mathews. ProbKnot: Fast prediction of RNA secondary structure including pseudoknots. *RNA*, 16(10):1870–1880, August 2010.
- [7] Steve L. Bonilla, Sarah K. Denny, John H. Shin, Aurora Alvarez-Buylla, William J. Greenleaf, and Daniel Herschlag. High-throughput dissection of the thermodynamic and conformational properties of a ubiquitous class of rna tertiary contact motifs. *Proceedings of the National Academy of Sciences*, 118(33):e2109085118, 2021.
- [8] Jessica A Brown, David Bulkley, Jimin Wang, Max L Valenstein, Therese A Yario, Thomas A Steitz, and Joan A Steitz. Structural insights into the stabilization of MALAT1 noncoding RNA by a bipartite triple helix. *Nature Structural & Molecular Biology*, 21(7):633–640, June 2014.
- [9] Nicola Calonaci, Alisha Jones, Francesca Cuturello, Michael Sattler, and Giovanni Bussi. Machine learning a model for RNA structure prediction. *NAR Genomics and Bioinformatics*, 2(4):1–10, 11 2020.
- [10] Xinshi Chen, Yu Li, Ramzan Umarov, Xin Gao, and Le Song. Rna secondary structure prediction by learning unrolled algorithms, 2020.
- [11] Isabel Chillón and Anna M. Pyle. Inverted repeat elements in the human lincRNA-p21 adopt a conserved secondary structure that regulates RNA function. *Nucleic Acids Research*, page gkw599, July 2016.
- [12] Uciel Chorostecki, Ester Saus, and Toni Gabaldón. Structural characterization of NORAD reveals a stabilizing role of spacers and two new repeat units. *Computational and Structural Biotechnology Journal*, 19:3245–3254, 2021.
- [13] Uciel Chorostecki, Jesse R. Willis, Ester Saus, and Toni Gabaldon. Profiling of RNA structure at single-nucleotide resolution using nextPARS. In *Methods in Molecular Biology*, pages 51–62. Springer US, 2021.
- [14] Padideh Danaee, Mason Rouches, Michelle Wiley, Dezhong Deng, Liang Huang, and David Hendrix. bpRNA: large-scale automated annotation and analysis of RNA secondary structure. *Nucleic Acids Research*, 46(11):5381–5394, May 2018.
- [15] K. E. Deigan, T. W. Li, D. H. Mathews, and K. M. Weeks. Accurate SHAPE-directed RNA structure determination. *Proceedings of the National Academy of Sciences*, 106(1):97–102, December 2008.
- [16] Janez Demsar. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research*, 7(1):1–30, 2006.
- [17] Yiliang Ding, Yin Tang, Chun Kit Kwok, Yu Zhang, Philip C Bevilacqua, and Sarah M Assmann. In vivo genome-wide profiling of rna secondary structure reveals novel regulatory features. *Nature*, 505(7485):696–700, 2014.
- [18] C.B. Do, D.A. Woods, and S. Batzoglou. CONTRAfold: RNA secondary structure prediction without physics-based models. *Bioinformatics*, 22(14):e90–e98, 2006.
- [19] Jörg Fallmann, Sebastian Will, Jan Engelhardt, Björn Grüning, Rolf Backofen, and Peter F. Stadler. Recent advances in RNA folding. *Journal of Biotechnology*, 261:97–104, November 2017.
- [20] Alessandro Fatica and Irene Bozzoni. Long non-coding RNAs: new players in cell differentiation and development. *Nature Reviews Genetics*, 15(1):7–21, December 2013.
- [21] Ryan A. Flynn and Howard Y. Chang. Long noncoding RNAs in cell-fate programming and reprogramming. *Cell Stem Cell*, 14(6):752–761, June 2014.

- [22] Laiyi Fu, Yingxin Cao, Jie Wu, Qinke Peng, Qing Nie, and Xiaohui Xie. Ufold: Fast and accurate RNA secondary structure prediction with deep learning. *Nucleic Acids Research*, 50(3):e14, 2022.
- [23] Boris Fürtig, Christian Richter, Jens Wöhnert, and Harald Schwalbe. NMR spectroscopy of RNA. *ChemBioChem*, 4(10):936–962, September 2003.
- [24] Johannes Graf and Markus Kretz. From structure to function: Route to understanding lncrna mechanism. *BioEssays*, 42(12):2000027, 2020.
- [25] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, November 1997.
- [26] Liang Huang, He Zhang, Dezhong Deng, Kai Zhao, Kaibo Liu, David A Hendrix, and David H Mathews. LinearFold: linear-time approximate RNA folding by 5′-to-3′ dynamic programming and beam search. *Bioinformatics*, 35(14):i295–i304, July 2019.
- [27] Ibrahim Avsar Ilik, Jeffrey J. Quinn, Plamen Georgiev, Filipe Tavares-Cadete, Daniel Maticzka, Sarah Toscano, Yue Wan, Robert C. Spitale, Nicholas Luscombe, Rolf Backofen, Howard Y. Chang, and Asifa Akhtar. Tandem stem-loops in roX RNAs act together to mediate x chromosome dosage compensation in drosophila. *Molecular Cell*, 51(2):156–173, July 2013.
- [28] S. Janssen and R. Giegerich. The RNA shapes studio. *Bioinformatics*, 31(3):423–425, October 2014.
- [29] Alisha N. Jones, Giuseppina Pisignano, Thomas Pavelitz, Jessica White, Martin Kinisu, Nicholas Forino, Dreycey Albin, and Gabriele Varani. An evolutionarily conserved RNA structure in the functional core of the lincRNA cyrano. *RNA*, 26(9):1234–1246, May 2020.
- [30] John Jumper, Richard Evans, Alexander Pritzel, and et al. Highly accurate protein structure prediction with AlphaFold. *Nature*, 596(7873):583–589, July 2021.
- [31] Yu-Jian Kang, De-Chang Yang, Lei Kong, Mei Hou, Yu-Qi Meng, Liping Wei, and Ge Gao. Cpc2: a fast and accurate coding potential calculator based on sequence intrinsic features. *Nucleic acids research*, 45(W1):W12–W16, 2017.
- [32] Philipp Kapranov, Jill Cheng, Sujit Dike, David A Nix, Radharani Dutttagupta, Aaron T Willingham, Peter F Stadler, Jana Hertel, Jorg Hackermuller, Ivo L Hofacker, et al. Rna maps reveal new rna classes and a possible function for pervasive transcription. *Science*, 316(5830):1484–1488, 2007.
- [33] Amanda Y. Keel, Robert P. Rambo, Robert T. Batey, and Jeffrey S. Kieft. A general strategy to solve the phase problem in RNA crystallography. *Structure*, 15(7):761–772, July 2007.
- [34] Michael Kertesz, Yue Wan, Elad Mazor, John L. Rinn, Robert C. Nutter, Howard Y. Chang, and Eran Segal. Genome-wide measurement of RNA secondary structure in yeast. *Nature*, 467(7311):103–107, September 2010.
- [35] Denise R Koessler, Debra J Knisley, Jeff Knisley, and Teresa Haynes. A predictive model for secondary RNA structure using graph theory and a neural network. *BMC Bioinformatics*, 11(S6):1–10, October 2010.
- [36] Soniya Lalwani and Rajesh Kumar. An efficient three-level parallel ABC algorithm for secondary structure prediction of complex RNA sequences. *Applied Soft Computing*, 99:106848, February 2021.
- [37] Hyunmin Lee, Zhaolei Zhang, and Henry M. Krause. Long noncoding RNAs and repetitive elements: Junk or intimate evolutionary partners? *Trends in Genetics*, 35(12):892–902, December 2019.
- [38] Jeannie T. Lee and Marisa S. Bartolomei. X-inactivation, imprinting, and long noncoding RNAs in health and disease. *Cell*, 152(6):1308–1323, March 2013.
- [39] Jung C. Lee and Robin R. Gutell. Diversity of base-pair conformations and their occurrence in rRNA structure and RNA structural motifs. *Journal of Molecular Biology*, 344(5):1225–1249, December 2004.
- [40] Fei Liu, Srinivas Somarowthu, and Anna Marie Pyle. Visualizing the secondary and tertiary architectural domains of lncRNA RepA. *Nature Chemical Biology*, 13(3):282–289, January 2017.
- [41] Zhendong Liu, Yurong Yang, Dongyan Li, Xinrong Lv, Xi Chen, and Qionghai Dai. Prediction of the rna tertiary structure based on a random sampling strategy and parallel mechanism. *Frontiers in genetics*, 12:813604, 2021.
- [42] Ronny Lorenz, Stephan H Bernhart, Christian Höner zu Siederdisen, Hakim Tafer, Christoph Flamm, Peter F Stadler, and Ivo L Hofacker. ViennaRNA package 2.0. *Algorithms for Molecular Biology*, 6(1):1–10, November 2011.

- [43] David Loughrey, Kyle E Watters, Alexander H Settle, and Julius B Lucks. Shape-seq 2.0: systematic optimization and extension of high-throughput chemical probing of rna secondary structure with next generation sequencing. *Nucleic acids research*, 42(21):e165–e165, 2014.
- [44] Weizhong Lu, Yan Cao, Hongjie Wu, Yijie Ding, Zhengwei Song, Yu Zhang, Qiming Fu, and Haiou Li. Research on RNA secondary structure predicting via bidirectional recurrent neural network. *BMC Bioinformatics*, 22(S3):1–10, May 2021.
- [45] Weizhong Lu, Ye Tang, Hongjie Wu, Hongmei Huang, Qiming Fu, Jing Qiu, and Haiou Li. Predicting RNA secondary structure via adaptive deep recurrent neural networks with energy-based filter. *BMC Bioinformatics*, 20(S25):1–10, December 2019.
- [46] Leandro Lucero, Lucía Ferrero, Camille Fonouni-Farde, and Federico Ariel. Functional classification of plant long noncoding RNAs: a transcript is known by the company it keeps. *New Phytologist*, 229(3):1251–1260, October 2020.
- [47] Kangkun Mao and Yi Xiao. Learning the fastest RNA folding path based on reinforcement learning and monte carlo tree search. *Molecules*, 26(15):4420, July 2021.
- [48] David H. Mathews. How to benchmark RNA secondary structure prediction accuracy. *Methods*, 162-163:60–67, June 2019.
- [49] David H. Mathews, Jeffrey Sabina, Michael Zuker, and Douglas H. Turner. Expanded sequence dependence of thermodynamic parameters improves prediction of RNA secondary structure. *Journal of Molecular Biology*, 288(5):911–940, May 1999.
- [50] Kevin V. Morris and John S. Mattick. The rise of regulatory RNA. *Nature Reviews Genetics*, 15(6):423–437, April 2014.
- [51] Stefanie A. Mortimer, Mary Anne Kidwell, and Jennifer A. Doudna. Insights into RNA structure and function from genome-wide studies. *Nature Reviews Genetics*, 15(7):469–479, May 2014.
- [52] V. L. Murthy. RNABase: an annotated database of RNA structures. *Nucleic Acids Research*, 31(1):502–504, January 2003.
- [53] Romasa Qasim, Nishat Kauser, and Tahseen Jilani. Secondary structure prediction of RNA using machine learning method. *International Journal of Computer Applications*, 10(6):15–22, November 2010.
- [54] Lijun Quan, Leixin Cai, Yu Chen, Jie Mei, Xiaoyu Sun, and Qiang Lyu. Developing parallel ant colonies filtered by deep learned constrains for predicting RNA secondary structure with pseudoknots. *Neurocomputing*, 384:104–114, April 2020.
- [55] Jessica S Reuter and David H Mathews. RNAstructure: software for RNA secondary structure prediction and analysis. *BMC Bioinformatics*, 11(1):1–10, March 2010.
- [56] John L. Rinn and Howard Y. Chang. Long noncoding RNAs: Molecular modalities to organismal functions. *Annual Review of Biochemistry*, 89(1):283–308, June 2020.
- [57] Gizem Rizki and Laurie A. Boyer. Lncding epigenetic control of transcription to cardiovascular development and disease. *Circulation Research*, 117(2):192–206, July 2015.
- [58] Natali Romero-Barrios, Maria Florencia Legascue, Moussa Benhamed, Federico Ariel, and Martin Crespi. Splicing regulation by long noncoding RNAs. *Nucleic Acids Research*, 46(5):2169–2184, February 2018.
- [59] Caroline Jane Ross and Igor Ulitsky. Discovering functional motifs in long noncoding RNAs. *WIREs RNA*, pages 1–10, January 2022.
- [60] K. Sato, M. Hamada, K. Asai, and T. Mituyama. CENTROIDFOLD: a web server for RNA secondary structure prediction. *Nucleic Acids Research*, 37(Web Server):W277–W280, May 2009.
- [61] Kengo Sato, Manato Akiyama, and Yasubumi Sakakibara. RNA secondary structure prediction using deep learning with thermodynamic integration. *Nature Communications*, 12(1):1–10, February 2021.
- [62] Kengo Sato and Yuki Kato. Prediction of RNA secondary structure including pseudoknots for long sequences. *Briefings in Bioinformatics*, 23(1):bbab395, 10 2022.
- [63] Kengo Sato, Yuki Kato, Michiaki Hamada, Tatsuya Akutsu, and Kiyoshi Asai. IPknot: fast and accurate prediction of RNA secondary structures with pseudoknots using integer programming. *Bioinformatics*, 27(13):i85–i93, June 2011.

- [64] Ester Saus, Jesse R. Willis, Leszek P. Pryszcz, Ahmed Hafez, Carlos Llorens, Heinz Himmelbauer, and Toni Gabaldón. nextPARS: parallel probing of RNA structures in illumina. *RNA*, 24(4):609–619, January 2018.
- [65] Susan J. Schroeder and Douglas H. Turner. Optical melting measurements of nucleic acid thermodynamics. In *Methods in Enzymology*, pages 371–387. Elsevier, 2009.
- [66] Matthew G. Seetin and David H. Mathews. RNA structure prediction: An overview of methods. In *Bacterial Regulatory RNA*, pages 99–122. Humana Press, 2012.
- [67] Jaswinder Singh, Jack Hanson, Kuldip Paliwal, and Yaoqi Zhou. RNA secondary structure prediction using an ensemble of two-dimensional deep neural networks and transfer learning. *Nature Communications*, 10(1):1–10, November 2019.
- [68] Jaswinder Singh, Kuldip Paliwal, Tongchuan Zhang, Jaspreet Singh, Thomas Litfin, and Yaoqi Zhou. Improved RNA secondary structure and tertiary base-pairing prediction using evolutionary profile, mutational coupling and two-dimensional transfer learning. *Bioinformatics*, 37(17):2589–2600, March 2021.
- [69] Matthew J. Smola, Thomas W. Christy, Kaoru Inoue, Cindo O. Nicholson, Matthew Friedersdorf, Jack D. Keene, David M. Lee, J. Mauro Calabrese, and Kevin M. Weeks. Shape reveals transcript-wide interactions, complex structural domains, and protein interactions across the xist lncrna in living cells. *PNAS*, 113(37):10322–10327, 2016.
- [70] Srinivas Somarowthu, Michal Legiewicz, Isabel Chillón, Marco Marcia, Fei Liu, and Anna Marie Pyle. HOTAIR forms an intricate and modular secondary structure. *Molecular Cell*, 58(2):353–361, April 2015.
- [71] Rachel Spokoini-Stern, Dimitar Stamo, Hadass Jessel, Lior Aharoni, Heiko Haschke, Jonathan Giron, Ron Unger, Eran Segal, Almogit Abu-Horowitz, and Ido Bachelet. Visualizing the structure and motion of the long noncoding rna hotair. *RNA*, 26(5):629–636, 2020.
- [72] Luisa Statello, Chun-Jie Guo, Ling-Ling Chen, and Maite Huarte. Gene regulation by long non-coding rnas and its biological functions. *Nature Reviews Molecular Cell Biology*, 22(2):96–118, 2021.
- [73] Zachary D. Stephens, Skylar Y. Lee, Faraz Faghri, Roy H. Campbell, Chengxiang Zhai, Miles J. Efron, Ravishankar Iyer, Michael C. Schatz, Saurabh Sinha, and Gene E. Robinson. Big data: Astronomical or genetical? *PLOS Biology*, 13(7):e1002195, July 2015.
- [74] Joanna Sztuba-Solinska, Jason W. Rausch, Rodman Smith, Jennifer T. Miller, Denise Whitby, and Stuart F.J. Le Grice. Kaposi sarcoma-associated herpesvirus polyadenylated nuclear RNA: a structural scaffold for nuclear, cytoplasmic and viral proteins. *Nucleic Acids Research*, 45(11):6805–6821, April 2017.
- [75] The ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature*, 489(7414):57–74, September 2012.
- [76] Ailone Tichon, Noa Gil, Yoav Lubelsky, Tal Havkin Solomon, Doron Lemze, Shalev Itzkovitz, Noam Stern-Ginossar, and Igor Ulitsky. A conserved abundant cytoplasmic long noncoding RNA modulates repression by pumilio proteins in human cells. *Nature Communications*, 7(1):1–10, July 2016.
- [77] Ailone Tichon, Rotem Ben-Tov Perry, Lovorka Stojic, and Igor Ulitsky. SAM68 is required for regulation of pumilio by the NORAD long noncoding RNA. *Genes & Development*, 32(1):70–78, January 2018.
- [78] Douglas H. Turner and David H. Mathews. NNDB: the nearest neighbor parameter database for predicting stability of nucleic acid secondary structure. *Nucleic Acids Research*, 38(suppl_1):D280–D282, October 2009.
- [79] Igor Ulitsky and David P. Bartel. lincRNAs: Genomics, evolution, and mechanisms. *Cell*, 154(1):26–46, July 2013.
- [80] Tina Uroda, Eleni Anastasakou, Annalisa Rossi, Jean-Marie Teulon, Jean-Luc Pellequer, Paolo Annibale, Ombeline Pessey, Alberto Inga, Isabel Chillón, and Marco Marcia. Conserved pseudoknots in lncRNA MEG3 are essential for stimulation of the p53 pathway. *Molecular Cell*, 75(5):982–995.e9, September 2019.
- [81] F. H. D. van Batenburg. PseudoBase: a database with RNA pseudoknots. *Nucleic Acids Research*, 28(1):201–204, January 2000.
- [82] Ashish Vaswani et al. Attention is all you need. In *Advances in Neural Information Processing Systems*, NIPS’17, page 60006010, Red Hook, NY, USA, 2017. Curran Associates Inc.

- [83] Linyu Wang, Yuaning Liu, Xiaodan Zhong, Haiming Liu, Chao Lu, Cong Li, and Hao Zhang. DMfold: A novel method to predict RNA secondary structure with pseudoknots based on deep learning and improved base pair maximization principle. *Frontiers in Genetics*, 10:1–10, March 2019.
- [84] Devin Willmott, David Murrugarra, and Qiang Ye. Improving RNA secondary structure prediction via state inference with deep recurrent neural networks. *Computational and Mathematical Biophysics*, 8(1):36–50, January 2020.
- [85] Ligang Wu and Joel G. Belasco. Let me count the ways: Mechanisms of gene regulation by miRNAs and siRNAs. *Molecular Cell*, 29(1):1–7, January 2008.
- [86] Haruka Yonemoto, Kiyoshi Asai, and Michiaki Hamada. A semi-supervised learning approach for RNA secondary structure prediction. *Computational Biology and Chemistry*, 57:72–79, August 2015.
- [87] Martina Zafferani and Amanda E. Hargrove. Small molecule targeting of biologically relevant rna tertiary and quaternary structures. *Cell Chemical Biology*, 28(5):594–609, 2021.
- [88] Shay Zakov, Yoav Goldberg, Michael Elhadad, and Michal Ziv-ukelson. Rich parameterization improves RNA structure prediction. *Journal of Computational Biology*, 18(11):1525–1542, November 2011.
- [89] Anna Zampetaki, Andreas Albrecht, and Kathleen Steinhofel. Long non-coding RNA structure and function: Is there a link? *Frontiers in Physiology*, 9:1–10, August 2018.
- [90] Hao Zhang, Chunhe Zhang, Zhi Li, Cong Li, Xu Wei, Borui Zhang, and Yuaning Liu. A new method of RNA secondary structure prediction based on convolutional neural network and dynamic programming. *Frontiers in Genetics*, 10:1–10, May 2019.
- [91] He Zhang, Liang Zhang, David H Mathews, and Liang Huang. LinearPartition: linear-time approximation of RNA folding partition function and base-pairing probabilities. *Bioinformatics*, 36(Supplement_1):i258–i267, July 2020.
- [92] Qi Zhao, Zheng Zhao, Xiaoya Fan, Zhengwei Yuan, Qian Mao, and Yudong Yao. Review of machine learning methods for RNA secondary structure prediction. *PLOS Computational Biology*, 17(8):e1009291, August 2021.
- [93] Omer Ziv, Svetlana Farberov, Jian You Lau, Eric Miska, Grzegorz Kudla, and Igor Ulitsky. Structural features within the NORAD long noncoding RNA underlie efficient repression of pumilio activity. *bioRxiv - Cold Spring Harbor Laboratory*, pages 1–20, November 2021.
- [94] Quan Zou, Jinjin Li, Qingqi Hong, Ziyu Lin, Yun Wu, Hua Shi, and Ying Ju. Prediction of MicroRNA-disease associations based on social network analysis methods. *BioMed Research International*, 2015:1–9, 2015.
- [95] 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.