

# Study of Accessible Motifs and RNA Folding Complexity

Yodo Wexler<sup>1,\*,\*\*</sup>, Chaya Zilberstein<sup>1,\*</sup>, and Michal Ziv-Ukelson<sup>2,\*</sup>

<sup>1</sup>School of Computer Science, Technion - Israel Institute of Technology, Haifa 32000, Israel

{ywex, chaya}@cs.technion.ac.il

<sup>2</sup>School of Computer Science, Tel-Aviv University, Tel-Aviv 69978, Israel

michaluz@post.tau.ac.il

**Abstract.** mRNA molecules are folded in the cells and therefore many of their substrings may actually be inaccessible to protein and microRNA binding. The need to apply an accessibility criterion to the task of genome-wide mRNA motif discovery raises the challenge of overcoming the core  $O(n^3)$  factor imposed by the time complexity of the currently best known algorithms for RNA secondary structure prediction [24, 25, 43].

We speed up the dynamic programming algorithms that are standard for RNA folding prediction. Our new approach significantly reduces the computations without sacrificing the optimality of the results, yielding an expected time complexity of  $O(n^2\psi(n))$ , where  $\psi(n)$  is shown to be constant on average under standard polymer folding models. Benchmark analysis confirms that in practice the runtime ratio between the previous approach and the new algorithm indeed grows linearly with increasing sequence size.

The fast new RNA folding algorithm is utilized for genome-wide discovery of accessible cis-regulatory motifs in data sets of ribosomal densities and decay rates of *S. cerevisiae* genes and to the mining of exposed binding sites of tissue-specific microRNAs in *A. Thaliana*.

Further details, including additional figures and proofs to all lemmas, can be found at: <http://www.cs.tau.ac.il/~michaluz/quadraticRNAFold.pdf>

## Introduction

The “lives” of messenger RNAs (mRNAs) begin with transcription and ultimately end with degradation. During their “lives”, mRNAs are translated into proteins. This process is regulated in a highly organized fashion to ensure that specific genes are expressed at the appropriate times and levels in response to various genetic and environmental stimuli [11, 35]. It is well-known that mRNA decay and translation are affected by regulatory motifs within mRNAs. These motifs serve as binding sites for transport proteins and microRNAs<sup>1</sup>. Several cis-regulatory RNA motifs were previously discovered experimentally, such as AREs (AU-Rich Elements) [28, 40], which

destabilizing elements involved in mRNA decay, and TOPs [13, 36], which inhibit the translation of ribosomal proteins and elongation factors.

Recently, new and interesting data has become available which measures, on a genome-wide scale, the ribosomal densities of mRNAs which reflect translation rates [37]. Additional data that measures mRNA decay rates [37]. The results of these measurements, if incorporated with genome-wide mRNA sequences, may reveal a wealth of cis-regulatory elements underlying both processes. However, since RNA elements are characterized by both *primary sequence* and higher order *structural conservation*, the identification of RNA elements is more complicated than identification of protein motifs. During the last decade, many computational efforts have been made to develop tools for the identification of RNA elements that are common to a group of orthologous or evolutionarily related genes. Some of these methods rely on a first step of multiple sequence alignment [2] and require that the sequences be highly similar globally, while other methods can detect locally conserved RNA sequence and structural elements in a subset of unaligned sequences [16, 26]. However, the complexity of these methods makes their application impractical for handling the large number of sequences involved in eukaryotic genome-wide analysis. Nevertheless, it turns out that most of the RNA regulatory motifs discovered so far are simple stem and loop structures, with a consensus motif residing in the loop area (e.g. IRES) [13, 36].

We note that the focus on *local* 2D structural conservation ignores the *global* accessibility of whether or not the primary sequence sites are indeed accessible to proteins. In order to allow the binding between the target cis-regulatory motif and trans-regulatory proteins or microRNAs, the base pairs in the motif must be free of any other chemical bond. This is due to the fact that the chemical recognition is based on the interaction between amino acids residing in the protein and the corresponding bases in the cis-regulatory motif residing in the mRNA [6], or on base pairing between the microRNA sequence and the motif nucleotides.

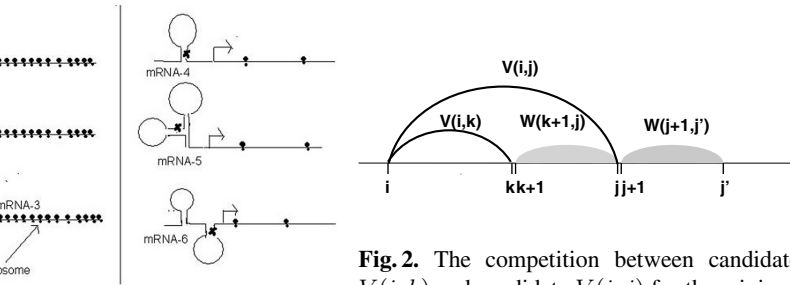
Our above requirement for chemical availability of motifs to protein binding calls for the formalization of an accessibility criterion:

**Definition 1 (“accessible” substring).** Let  $S$  be a sequence and  $s$  a region i.e. substring of  $S$ . We say that  $s$  is **accessible** iff the following two conditions apply:

1. There exists a 2D structure of  $S$  with predicted free energy  $G_1$  in which none of the nucleotides of  $s$  is engaged in base pairing.

2.  $G_1 - G_0 \leq \delta$ , where  $\delta$  is a user defined threshold parameter, and  $G_0$  is the optimal free energy of the full string  $S$ .

Here we suggest a novel approach to the genome-wide discovery of RNA cis-regulatory motifs. In our framework, motifs are scored according to their statistical significance when applying the above accessibility criterion. In order to accommodate high-throughput mRNA sequences are first filtered according to Definition 1. This is done to reduce the noise created by motifs which are not exposed to trans-regulatory



**Fig. 2.** The competition between candidate  $V(i, k)$  and candidate  $V(i, j)$  for the minimal  $W(i, j')$ . Candidate  $V(i, k)$  has an advantage over candidate  $V(i, j)$  in the additional potential cost for segment  $s_{j+1} \dots s_{j'}$  since it has a wider left-scope for combining this segment in a structure with  $W(k+1, j)$ . Therefore, if  $V(i, k) + W(k+1, j) \leq V(i, j)$  then by triangle inequality  $V(i, k) + W(k+1, j') \leq V(i, j) + W(j+1, j')$ .

Applying the accessibility criterion to aid motif discovery in mRNA density data. The motif  $X$  may be used to differentiate between the set of motifs with high density (left) and the set of motifs with low density (right) since its occurrences in motifs 1, 4, 5 and 6 are inaccessible.

Next, the mRNA corresponding to the gene which is targeted for “knock out”, is scanned for accessible sites. For this task, the current RNA folding algorithms are sufficient. However, such tools could not be practically scaled up to whole genome motif discovery, where thousands of mRNAs need to be mined for accessible sites, without raising severe efficiency problems: the complexity of RNA folding prediction allowing multiple loops but no pseudoknots is  $O(n^3)$  to begin with, where  $n$  is the size of an RNA sequence (typically  $\sim 2000$ ). This complexity is further increased to  $O(n^3 \cdot m)$  by the need to exhaustively run a sliding window across the genome, where  $m = O(n)$  is the number of different starting positions of accessible regions that need to be considered in each gene. Note that the sliding window challenge is not addressed by Robins et al. [27], where the computation is complicated by the fact that only a single optimal folding is computed per gene. Thus, the task of mining accessible sites for genome-wide motif discovery creates a heavy computational bottleneck in terms of computational complexity, where  $g$  is the number of genes in the genome under study (typically in the thousands).

Practical considerations raised by such a complexity are exemplified as follows: if the genome under study contains 6000 mRNA sequences, of size  $\sim 2000$  nucleotides each, in which we need to consider all potential sites obtained by sliding a window of size  $k \ll 2000$ . Given that the folding prediction computation for each window takes about twenty seconds<sup>2</sup>: the total time needed for the computation of all accessible sites in this case would be  $6000 \cdot 2000 \cdot 20$  seconds  $\approx 7.61$  years! Note that even if we confine our search to  $\sim 300$  windows in the UTR regions, the time needed still sums up to more than a year. This example demonstrates the need for efficient folding algorithms, especially when dealing with whole-genome scale data.

Wexler, C. Zilberstein, and M. Ziv-Ukelson

classical  $O(n^3)$  algorithms for RNA secondary structure prediction [25, 43], have been heavily used by the bioinformatics community in the last two decades, substantially sped up? Furthermore, could such a speed up be implemented via a linear, low-constant algorithm?

Important challenges are addressed in the rest of this paper, where we describe a dynamic programming algorithm that exploits the combination of two properties of RNA secondary structure prediction: one is the observed *triangle inequality* of the matrices commonly used in RNA secondary structure prediction (Section 2.1) and the other is the *polymer-zeta* behavior of RNA folding with respect to sequence size (Section 2.4). These observations are utilized here via a simple list algorithm, called Algorithm CANDIDATEFOLD (Section 2.3), which significantly reduces the computations without sacrificing the optimality of the results (no approximations are used). The expected time complexity of Algorithm CANDIDATEFOLD is  $O(\psi(n))$  instead of the previously known  $O(n^3)$ , where  $\psi(n)$  is shown to converge to  $O(n)$  under models previously described for RNA folding and re-validated by experiments (see Section 2.5). Furthermore, due to the simplicity of Algorithm CANDIDATEFOLD, it is indeed much faster than the classical algorithm in practice, as supported by experimental performance results in Section 3. Clearly, this new algorithm for RNA folding prediction is applicable to a wide range of additional biological applications, especially to those that require a substantial amount of RNA folding computations.

On the efficient new RNA folding algorithm CANDIDATEFOLD, we conducted an analysis which examines the contribution of the “accessible site” criterion to the discovery of RNA motifs that would otherwise be obscured by noise. The new approach was applied to quantitative data sets of ribosomal densities and decay rates of almost all 1000 *S. cerevisiae* genes. By applying our approach, some biologically interesting motifs were discovered (Section 5). For example, the motif *AGCKTTA* in the decay rates data was  $5 \cdot 10^{-7}$ . This  $p$ -value was significant because the fact that the average half-life (i.e.  $\log(2)/\text{decay rate}$ ) of 24 genes that were found to contain this motif in an accessible substring was 26 days, while the half-life of the background population was 15 days. Relaxing the accessibility criterion lowered the significance of the motif by raising its  $p$ -value to 0.008.

We also employed the “accessible target” criterion to analyze microRNAs regulating biological processes in *A. Thaliana*. Interesting tissue specific microRNAs were identified (see Fig. 4).

## Accessible Site Prediction Engine

### Foundations of RNA Folding Prediction Via Minimum Energy

RNA is typically produced as a single stranded molecule which then folds intracellularly to form a number of short base-paired stems. This base-paired structure is the thermodynamic minimum energy state of the RNA molecule.

s which are standard for RNA structure prediction do not deal with pseudoknots. This is done mostly in order to simplify the problem and is justified by the fact that pseudoknots do not contribute much to the overall energy and long pseudoknots are kinetically difficult to form [20]. Therefore, in this paper we assume that no pseudoknots are allowed, however multiple loops are indeed allowed.

Under the above assumptions, a model was proposed in Tinoco et al. [32] to calculate the thermodynamic stability (in terms of free energy) of a folded RNA molecule by adding contributions from base pair stacking and loop-destabilizing terms from the RNA secondary structure. This model has proven to be a good approximation of the forces governing RNA structure formation, thus allowing fair predictions of real structures by computing the most stable structures in the model of a given sequence. Based on this model, several algorithms for computing the most stable structures have been proposed (Nussli and Jacobson, 1980 [25]; Zuker and Steigler, 1981 [43]), and various tools for RNA secondary structure prediction were developed. The tools commonly used today are RNAfold [42], Vienna Package [14] and FOLD RNA [41].

The thermodynamic parameters used by our accessible site prediction engine are the same as those used by the RNA folding tools listed above. In this paper, the following four recursions are combined to model RNA secondary structure folding. Note that the recursions depend on the nature of the energy rules used. Where  $eh(i, j)$  is the energy of the hairpin loop closed by the base pair  $i, j$ ,  $es(i, j)$  is the energy of the stacked pair  $i, j$  and  $i + 1, j - 1$  and  $ebi(i, j, i', j')$  is the energy of a bulge or an interior loop closed by  $i, j$  with  $i', j'$  accessible from  $i, j$ . The boundary conditions  $W(i, j) = V(i, j) = +\infty$  if  $j - i < 4$ . More recursions, based on the ones given here, take into consideration exterior base pairs [43]. These are not elaborated here for the sake of simplicity of presentation, but the same reasoning applies to this extension as well. The recursion equations are stated below:

$$W(i, j) = \min\{V(i, j), W(i + 1, j), W(i, j - 1), \min_{i \leq k < j} \{W(i, k) + W(k + 1, j)\}\} \quad (1)$$

Equation (1) computes the optimal folding of substring  $s_i, \dots, s_j$ , which is the value of the entry at row  $i$  and column  $j$  of the main, upper-triangular DP table  $W$ . The computation of  $W$  involves the matrix  $V$  whose entries are computed via the following recursion:

$$V(i, j) = \min\{eh(i, j), es(i, j) + V(i + 1, j - 1), VBI(i, j), VM(i, j)\} \quad (2)$$

Equation (2) computes the optimal folding energy of a substring  $s_i \dots s_j$  in which  $s_i$  base pairs with  $s_j$ .

$$VBI(i, j) = \min_{i < i' < j' < j} \{ebi(i, j, i', j') + V(i, j)\} \quad (3)$$

Equation (3) computes the score of an optimal folding of substring  $s_i, \dots, s_j$  given that there is a base pair formed at indices  $(i, i', j', j)$ .

**Analysis of the Classical RNA Folding Prediction Engine.** The above recursion is implemented by maintaining four tables of size  $O(n^2)$  each. Eq. 1 is clearly linear in the values computed for Eq. 1, the values for Eq. 4 can be computed in constant time and space via direct look-up of the minima values previously computed for Eq. 2 is also  $O(n^2)$ .

The complexity for the computation of internal loop size energies is naively  $O(n^4)$ . Practically, it is standard to assume that RNA interior loop size is bounded by a constant (15 nt at room temperature and up to 30 nt in extreme heat). The program RNAFOLD in the Vienna package [14] as well as the MFOLD program [42] use constant gap size in both directions to reduce the complexity of Eq. 3 to  $O(n^2)$ . Lygnso *et. al.* [22] show how to reduce the complexity of this equation to  $O(n^3)$  without binding the gap size. On the theoretical front, Waterman and Smith showed how to compute internal loops in  $O(n^3)$  assuming that the loop penalty is a function of its size [34]. Eppstein, Galil and Wigderson [7, 9] considered loop destabilizing functions satisfying certain convexity or concavity conditions, and developed an  $O(n^2 \log^2 n)$  algorithm for this case. This was improved to  $O(n^2 \log n)$  [1], and finally to  $O(n^2 \alpha(n))$  (where  $\alpha$  is the inverse of Ackermann's function) for logarithmically growing destabilizing functions [19].

**Lemma 1.** *The  $O(n^3)$  bottleneck to RNA Folding Prediction complexity is based on the minimization term  $\min_{i \leq k < j} \{W(i, k) + W(k + 1, j)\}$  in Eq. 1.*

The  $O(n^3)$  bound applies to both the *worst case* and the *expected case* time complexities of the classical RNA folding algorithm, since Eq. 1 is called  $O(n^2)$  times and each call involves the computation of the minimum over  $O(n)$  elements on average.

## Triangle Inequality in the Context of Dynamic Programming

In this section we formalize the *triangle inequality* property in the context of dynamic programming tables and show that the main matrix  $W$ , which is the final output of the dynamic programming recursions given in the previous section, obeys this property. Let  $M$  be a matrix in which each entry  $M(i, j)$  ( $i \leq j$ ) is computed by the following formula:

$$M(i, j) = \min_{i < k \leq j} \{M(i, k) + M(k + 1, j)\}$$

The known inverse quadrangle inequality property [10] is defined as follows.

**Lemma 2.** *A matrix  $M$  obeys the inverse quadrangle inequality condition iff*

$$i < i' < j < j' \quad M(i, j') \leq M(i, j) + M(i', j') - M(j', j)$$

The quadrangle and the inverse quadrangle inequalities have previously been used in the context of dynamic programming [5, 10]. However, both the quadrangle inequality and the inverse quadrangle inequality are strong constraints on the input behavior, and do not hold for arbitrary RNA folding energy functions. For example, the energy function

**3.** A matrix  $M$  obeys the **triangle inequality property** iff

$$\forall i < j < j' \quad M(i, j') \leq M(i, j) + M(j + 1, j').$$

### Simple 1D Candidate List Approach to the Construction of $W$

$s_1 \dots s_n$  denote a given RNA sequence. The next two definitions describe folding concepts that will be used in the description of the new algorithm.

**4 (Structure).** A **structure** over a sequence  $s_i \dots s_j$  is a folding in which pairs with  $s_j$ .

**5 (Partition Point).** A **partition point** in a given folding of  $S = s_1 \dots s_n$  is, such that there is no structure over  $s_i \dots s_j$  in this folding, where  $1 \leq i \leq k \leq n$ .

In this section we describe an alternative approach to the computation of  $W$ , which is simpler than the standard algorithm. Similarly to the standard algorithm, the new algorithm computes the value of  $W(i, j)$  row by row, in bottom-up order (decreasing row index). For each row  $i$  of  $W$ ,  $W(i, j)$  is computed in left-to-right order (increasing column index). However, the suggested new algorithm, called CANDIDATEFOLD, differs from the original algorithm in the application of Eq. 1 to the computation of  $W(i, j)$ . In a given row  $i$ , instead of considering  $O(n)$  possible partition points for each column  $j$  in Eq. 1, the new algorithm only considers a list of candidate partition points, which are maintained in the form of a simple candidate list. In the following sections we show that the expected size of this candidate list for an  $n$ -sized sequence, denoted  $\psi(n)$ , is constant. To more clearly define the properties that make a potential partition point a qualified candidate, we first need to simplify Eq. 1. Note that, if the main diagonal  $W(r, r)$  is zero, then the two terms  $W(i + 1, j)$  and  $W(i, j - 1)$  in Eq. 1 could be moved into the minimization term as special cases.  $W(i + 1, j)$  would then be obtained as the special case  $k = i + 1$  to yield the sum  $W(i, i) + W(i + 1, j)$  which is exactly  $W(i, j)$ ; similarly,  $W(i, j - 1)$  would be obtained as the special case  $k = j - 1$  to yield the sum  $W(i, j - 1) + W(j, j)$  which is exactly  $W(i, j - 1)$ . However, the problem is that setting  $W(r, r) = 0$  would contradict the boundary conditions set by Stiegler [43], which assume that  $W(r, r) = \infty$ .

Therefore, we add two auxiliary matrices, denoted  $W'$  and  $V'$ , computed via the recurrence as given below, where Eq. 7 replaces the previous Eq. 1. Note that the matrix  $W'$  is used in order to get around the above boundary condition problem, while matrix  $V'$  is used to simplify the presentation of the algorithm which is described in the next

$$W(i, j) = W'(i, j) \quad \forall j \geq i + 4 \quad (5)$$

$$V'(i, j) = V(i, j) \quad \forall j \geq i + 4 \quad (6)$$

$$W'(i, j) = \min\{V'(i, j), \min_{i \leq k < j} \{W'(i, k) + W'(k + 1, j)\}\} \quad (7)$$

the values of  $W(i, j)$  and  $V(i, j)$ , as computed via Eqs. 2-7, are identical to those obtained when using Eqs. 1-4.

This claim is immediate from Definition 2 and Eq. 7.

The matrix  $W'$ , as computed by Eq. 7, obeys the triangle inequality.

This claim is used in the next lemma to show that any sum which yields the optimal score of Eq. 7 can be reformulated as a corresponding, equal-scoring sum, in which the intermediate structure is a structure (see Definition 4).

**Lemma 1.** *Consider Eq. 7. For every entry  $W'(i, j)$ , if there exists an index  $k$ ,  $i \leq k < j$ , such that  $V'(i, j) = W'(i, k) + W'(k + 1, j)$ , then  $W'(i, k') = V'(i, k')$  for some index  $k' < j$ .*

Using Lemma 1, Eq. 7 can be reformulated as follows.

$$W'(i, j) = \min\{V'(i, j), \min_{i \leq k < j} \{V'(i, k) + W'(k + 1, j)\}\} \quad (8)$$

However, after the transformation to Eq. 8, there are still  $n$  candidate partition points that must be compared to compute the optimal score in the minimization term. However, the next theorem establishes a dominance relationship between these candidates (see Figure 2).

**Lemma 2.** *If  $V'(i, j) \geq V'(i, k) + W'(k + 1, j)$  for some  $i < k < j$ . Then,*

$$j' > j \implies V'(i, j) + W'(j + 1, j') \geq V'(i, k) + W'(k + 1, j').$$

Lemma 2 exposes redundancies in the  $O(n)$  computation of Eq. 8, which could be avoided by maintaining a list of only those candidates that are not dominated by others.

**Definition 6 (candidate).** *A column index  $j$  is a **candidate** in a row  $i \leq j$  iff  $V'(i, j) < V'(i, k) + W'(k + 1, j) \forall i \leq k < j$ .*

This definition can be applied to speed up the computation of  $W'(i, j)$ , as follows. Instead of considering all possible  $n$  partition point indices for the computation of  $W'(i, j)$ , one could query the list that contains only partition points that satisfy the candidate criterion according to Definition 6. This is formalized in the following equation,

$$W'(i, j) = \min\{V'(i, j), \min_{\forall k \in \text{candidate\_list}} \{V'(i, k) + W'(k + 1, j)\}\} \quad (9)$$

This is implemented via a candidate list that is empty at the start of each row and is updated throughout the left-to-right computation of row  $i$  by appending only those partition points which are candidates by Definition 6. Each partition point is considered a candidate once per row, when its column is reached. The pseudo-code for the algorithm for computing Eq. 7, denoted *Algorithm CANDIDATEFOLD*, is given below.



$$W'(i, j) \leftarrow \min_{k \in \text{candidate\_list}} \{ V'(i, k) + W'(k + 1, j) \}$$

if ( $V'(i, j) < W'(i, j)$ ) then  
 $W'(i, j) \leftarrow V'(i, j)$   
Append  $j$  to the *candidate\_list*

### Case Time Analysis of the Improved RNA Folding Prediction Engine.

denote the expected maximal size of the candidate list in a sequence of size  $n$ . CANDIDATEFOLD computes each entry in the  $n^2$ -sized energy-matrix  $W'$ . Each calculation requires the computation of Eq. 9, where the major work is that finding the minimum among  $O(\psi(n))$  candidates. All other recursions remain constant. Therefore, the overall average time complexity is  $O(n^2 \cdot \psi(n))$  if the standard on interior loop size is followed, or otherwise  $O(n^2 \cdot \max\{\psi(n), \alpha(n)\})$ , where  $\alpha$  is the inverse ackerman function.

In next sections we analyze the expected growth of the candidate list size with increasing sequence size and assert the surprising fact that  $\psi(n)$  converges to a constant. This leads to the conclusion that Algorithm CANDIDATEFOLD improves over  $O(n^3)$  classical algorithm (analyzed in section 2.1) by a linear factor on

### Polymer-Zeta Property of RNA Folding

The polymer-zeta property is defined as follows.

**Definition 7.** Let  $P(i, j)$  denote the probability of a structure over the substring  $i \dots j$  under a given set  $\Lambda$  of folding rules, where  $j - i = m$ . We say that  $\Lambda$  follows the **polymer-zeta property** if  $P(i, j) = b/m^c$  for some constants  $b, c > 0$ .

Our work shows that RNA, which folds like other polymers, obeys the polymer-zeta property, namely, the probability that a structure is formed over the subsequence between two positions distant  $m$  monomers apart is  $P(m) = b/m^c$  where  $b = 1$  and  $c = 1.5$  [18]. This fact is explained by modeling the 2D folding of a polymer chain as a self avoiding random walk (SAW) in a 2D lattice [33]. In this model the spacial position of every nucleotide in the original polymer corresponds to a random step in the lattice, where edges of the lattice represent possible transition directions. Since this model of polymer folding also ignores pseudoknots, the walk is called “self avoiding”, and the assumption is followed that the walk does not intersect the prefix of the chain. The interest here is the probability that the  $m^{th}$  step in the self avoiding random walk occupies the same vertex in the lattice as the origin. The theoretical exponent for the 2D SAW model is known to be  $c = 1.5$  [8]. This is supported in Monte Carlo simulations for collapsing polymers of sequence size up to 3200, as reported in [18]. These simulations exhibited an exponent of 1.375 at low temperatures and 1.571 at high temperatures.

The thermodynamic programming algorithm follows the thermodynamic rules defined by

Wexler, C. Zilberstein, and M. Ziv-Ukelson

single structure formation probabilities in polymer folding, which were found to be a polymer-zeta property. We used 50,000 mRNA sequences with an average length of 992 nucleotides from the NCBI databases and found that the probability that a random folding forms a structure over  $s_i \dots s_j$ , where  $m = j - i$ , is estimated to be  $b \cdot m^{-c}$ . The degree exponent  $c$  was estimated in our study to be  $\sim 1.47$  by standard statistical procedures (approximating the MLE parameter followed by “Kolmogorov-Smirnov” and “chi-square” goodness-of-fit tests, using the *R* analysis package, <http://www.r-project.org>).

## Conclusions on $\psi(n)$

We analyze  $\psi(n)$  based on our findings. The following observation is immediate (Lemma 1).

**Lemma 1.** *A new candidate  $j$  is added to the candidate list, in step 6 of Algorithm CANDIDATEFOLD, iff the optimal predicted folding of substring  $s_i \dots s_j$  forms a simple structure from index  $i$  to index  $j$ . The only exception to this case is the boundary candidate  $i$ , which is always added as a “virtual” structure to the list.*

Let the probability for a new candidate situated  $m$  bases away from the start of a sequence is  $b \cdot m^{-c}$ , the expected number of candidates in a sequence of length  $n$  is  $b \sum_{i=1}^n i^{-c}$ . This summation could assume one of three values, according to the value of  $c$ :

If  $c < 1$ , this series is a partial sum of the *Riemann Zeta function* defined as  $\sum_{i=1}^{\infty} i^{-c}$ .

If  $c > 1$ , this series is known to converge and thus,  $\psi(n) = O(1)$ .

If  $c = 1$ , we get a partial sum of the first  $n$  elements of the Harmonic series, which is known to be less or equal to  $1 + \ln(n)$  and thus  $\psi(n) = O(\log n)$ .

If  $c > 1$ , we use the power means inequality to obtain the bound  $\psi(n) = O(n^{1-c})$ .

**Lemma 2.** *Applying Algorithm CANDIDATEFOLD to the folding of a polymer chain that obeys the polymer-zeta property with  $c > 1$ , requires an average of  $O(n^2)$  operations.*

Our simulations estimate  $c$  to be 1.47, which implies that  $\psi(n) \sim 2.11 \cdot n^{-0.47}$ , which is a constant. Therefore, applying Algorithm CANDIDATEFOLD to the folding of an RNA sequence of size  $n$  takes  $O(n^2)$  time on average.

## Performance of the New RNA Folding Engine

To demonstrate the power of algorithm CANDIDATEFOLD in practice we ran it against a set of RNA sequences of varying lengths and compared its performance to that of the standard RNA folding algorithm.

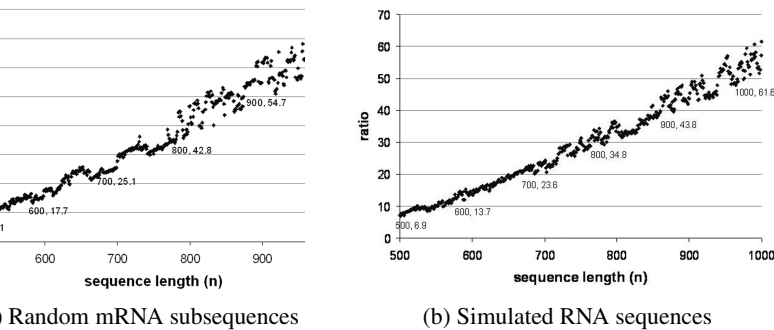


Figure 3. Average measured run-time ratio of naive/CANDIDATEFOLD as a function of increasing sequence size

Figure 3 demonstrates that the average run time ratio (computed by dividing the run time of the classical algorithm with ours) is linear in the sequence length  $n$ , reinforcing our time complexity analysis. In Figure 3(a), the analysis was done for 100 sequences of each size in the range 500-1000, which were extracted as random subsequences from 50,000 complete mRNA sequences taken from NCBI. The analysis shown in Figure 3(b) was done for 100 sequences of each size in the same range, which were generated using a Markov-model imitating software. The sequence-simulation program takes a set of sequences to imitate and a Markovian order parameter, and generates an output of random sequences according to a Markovian process of the desired order. The input consisted of 50,000 complete mRNA sequences extracted from the NCBI database and the Markovian order parameter was set to 6. Similar results emerged when using the remaining 50,000 mRNA sequences as input to the order Markovian model simulator.

## Methods for Mining Accessible Cis and Trans Regulatory Motifs

Our method for discovering novel cis-regulatory motifs incorporates large scale decay rate and ribosomal density measurements, combined with the information from mRNA accessibility. The method can be formulated as follows. Given a set of sequences  $G = S_1 \dots S_g$ , a parameter  $k$  denoting motif window size (could be slightly larger than the motif residing in the window), and a pre-defined energy threshold  $\delta$ , we follow the following simple two-stage approach:

1. Process the sequence set  $G$  to extract all “accessible” windows by running a sliding window of size  $k$  across the mRNA sequence and testing each window for compliance with Definition 1. For each shifted window this testing is conducted by masking nucleotides inside the window in order to prevent their engagement in base pairing.

Wexler, C. Zilberstein, and M. Ziv-Ukelson

This stage takes as input the accessible substrings, extracted in the first stage, regulatory motifs residing in the data. Two statistical techniques are applied depending on whether the sought motif is cis or trans regulatory:

**Transitory motifs:** Enumerate all motifs up to a given size  $k$  over the IUPAC alphabet [8]. For each motif use the new data created in stage 1 instead of the original sequences, to compute a  $t$ -score [12] reflecting the functionality of that motif. If the  $t$ -score associated with the computed  $t$ -score is small enough, report the motif. This can be efficiently executed by using a variation of the algorithm of Sagot combined with the statistical computation of the  $t$ -score [38] and adapted to the new “accessible window” data.

**Regulatory Signals (microRNAs).** The search for microRNAs is similar to that of motifs, except for the following difference: instead of considering accessible motifs, we considered accessible sites that were predicted to hybridize well with microRNAs, as described in [39].

## Biological Study of Accessible Regulatory RNA Elements

We conducted a study in order to test our novel approach, which applies the “accessibility” criterion to RNA motif discovery. Using various data sets, significant motifs were discovered, including some cis-regulatory degradation and translation motifs and specific microRNAs.

Of the conducted experiments, two data sets were studied: a set containing “accessible” substrings, according to Definition 1, and a “control” set which includes the original complete mRNA sequences. A comparison of the results obtained from the two sets repeatedly confirms the contribution of the “accessibility” criterion as a powerful filter for masking out noise associated with inaccessible motifs and for increasing the significance score of otherwise invisible motifs.

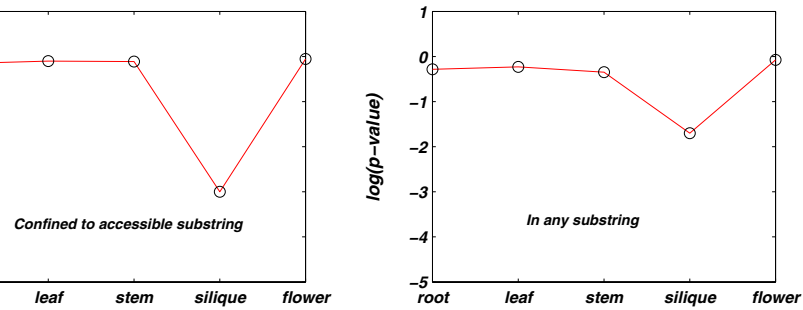
**On Related Motifs.** Arava *et al.* [3] measured the ribosomal densities of all mRNAs of the yeast *S. cerevisiae* under normal cell conditions, using the sucrose gradient method. First, mRNAs are extracted from the cells and separated by velocity sedimentation. Then, each fraction across the gradient is analyzed by microarray technology for its mRNA content. Based on this, a fraction is assigned to each mRNA: the higher the fraction is, the higher the mRNA’s ribosomal density is. We applied our approach to this data in order to detect translation cis-regulatory elements within 5’ UTR region (5’UTR)<sup>3</sup>. A few novel potential cis-regulatory elements were discovered that may affect translational efficiency (see Table 1). In particular, the average ribosomal density of the set of mRNAs containing the motif *AGSNNK* in accessible regions was low in comparison to the background. Thus, *AGSNNK* seems to be a repressor.

motifs potentially regulating mRNA translations. The accessible substring criterion was with window size 10 and  $\delta = 2Kcal$ . The average ribosomal density without the motif was computed based on  $\sim 5000$  different genes.

	Number of occurrences	Average density with the motif	Average density without the motif	p-value confined to accessible substrings	p-value in any substring	Hypothesized function
TT	14	1.7	0.7	$10^{-18}$	$10^{-4}$	Translation enhancer
KK	1292	0.6	0.7	$10^{-11}$	$10^{-3}$	Translation repressor

motifs potentially regulating mRNA degradations. The first 3 columns refer to the case of the substring with window size 10 and  $\delta = 2Kcal$ . The average half life without the motif was computed based on  $\sim 5000$  different genes.

	Number of occurrences	Average half-life with the motif	Average half-life without the motif	p-value confined to accessible substrings	p-value in any substring	Hypothesized function
AA	24	26.54	15.46	$4.83 \cdot 10^{-7}$	0.0083	Stabilizer
RR	5	57.75	15.5	$2.76 \cdot 10^{-9}$	0.0081	Stabilizer
TT	4	42.75	15.49	$4.84 \cdot 10^{-7}$	0.01198	Stabilizer



-161 and it's p-values in different plant tissues. The accessible substring criterion was with window size 25 and  $\delta = 6Kcal$ .

regulating elements within 3' UTRs<sup>4</sup>. We successfully identified some novel cis-regulatory motifs that may affect mRNA stability (see Table 2). For example, the average half-lives (i.e.  $\log(2)/\text{Decay rate}$ ) of the set of mRNAs containing the motif *AGCKTTA* in accessible substrings was high in comparison to the control. Thus, *AGCKTTA* seems to be a strong mRNA stabilizer. Table 2 also indicates that, when relieving the accessibility criterion, the significance of the p-value substantially dropped.

**Specific microRNAs.** In order to discover microRNAs, which are potential trans-regulating mRNA stabilities, we collected the genome-wide expression

Wexler, C. Zilberstein, and M. Ziv-Ukelson

ere discovered<sup>5</sup>. These microRNAs showed a significant  $p$ -value for binding the tissues and non-significant  $p$ -values in the rest of the tissues. For example, RNA *miR*-161, represented in Figure 4, is specific to silique tissue. Interest-figure demonstrates that in most of the tissues the  $p$  – values corresponding t (accessible substring) and second (control) input sets are almost similar. in the silique tissue, where the microRNA *miR*-161 seems to be active, the between the two input sets becomes conspicuous.

**edgments.** Many thanks to Yoav Arava for inspiration and data, as well as l discussions. We thank Micheal Zuker for very helpful advice. The authors grateful to Ron Shamir, Ron Y. Pinter, Dan Geiger, Zohar Yakhini, Jeannette Christos Faloutsos, Eleazar Eskin and Firas Swidan for helpful discussions ents. The research of Michal Ziv-Ukelson was supported in part by the Aly Post Doctoral Fellowship.

## ces

garawal and J. Park. Notes on searching in multidimensional monotone arrays. *Proc. IEEE Symp. on Foundations of Computer Science*, 497–512, 1988.

maev, S. Kelley, and G. Stormo. A phylogenetic approach to RNA structure prediction. *Int Conf Intell Syst Mol Biol*, 235:10–17, 1999.

ya, Y. Wang, J. Storey, C. Liu, P. Brown, and D. Herschlag. Genome-wide analysis of translation profiles in *saccharomyces cerevisiae*. *PNAS*, 100:3889–3894, 2003.

istofferson et al. Application of computational technologies to ribozyme biotechnol- products. *J.Molecular Struct.(Theochem)*, 311:273, 1994.

chemore, G. Landau, B. Schieber, and M. Ziv-Ukelson. *Re-Use Dynamic Program- or Sequence Alignment:An Algorithmic Toolkit. String Algorithmics*. KCL Press, 2005.

per. Themes in RNA-protein recognition. *J Mol Biol*, 293(2):255–270, 1999.

stein, Z. Galil, and R. Giancarlo. Speeding up dynamic programming. *Proc. 29th Symp. on Foundations of Computer Science*, 488–496, 1988.

ner. Shape of a self-avoiding walk or polymer chain. *JCP*, 44:616–622, 1966.

l and R. Giancarlo. Speeding up dynamic programming with applications to molecular y. *Theoretical Computer Science*, 64:107–118, 1989.

ncarlo. *Dynamic Programming: Special Cases*. In *Pattern Matching Algorithms*, A. Gallo and Z. Galil eds., Oxford University Press, 1997.

win, P. Okkema, T. C. Evans, and J. Kimble. Translational regulation of tra-2 by its translated region controls sexual identity in *c. elegans*. *Cell*, 75:329–339, 1993.

lden. *Methods of Statistical Analysis*. New York: Wiley, 2 edition, 1956.

y and M. Wickens. *Annu Rev Cell Dev Biol*, 14:399–458, 1998.

ofacker. Vienna RNA secondary structure server. *NAR*, (13):3429–3431, 2003.

araman and S.P.Walton. Rational selection and quantitative evaluation of antisense nucleotides. *Biochim.Biophys. Acta*, 1520:105, 2001.

X. Xu, and G. Stormo. *Bioinformatics*, 20:1591–1602, 2004.

akcioglu and A. Stella. A scale-free network hidden in the collapsing polymer. *ArXiv preprint*, 2004.

- ari, D. Mukamel, and L. Peliti. Why is the dna denaturation transition first order? *Physical Review Letters*, 85:4988–4991, 2000.
- more and B. Schieber. On-line dynamic programming with applications to the prediction of RNA secondary structure. *J. Algorithms*, 12(3):490–515, 1991.
- and R. Bundschuh. Quantification of the differences between quenched and annealed RNA secondary structures. *ArXiv Physics e-prints*, Apr. 2005.
- ve et al. Cleavage of scarecrow-like mRNA targets directed by a class of arabidopsis proteins. *Science*, 297:2053–2056, 2002.
- lyngsø, M. Zuker, and C. N. S. Pedersen. An improved algorithm for RNA secondary structure prediction. Technical Report RS-99-15, brics, 1999.
- hews et al. *RNA*, 5, 1458–1469, 1999.
- hews, J. Sabina, M. Zuker, and D. Turner. *JMB*, 288:911, 1999.
- sinov and A. Jacobson. Fast algorithm for predicting the secondary structure of single-stranded RNA. *PNAS*, 77(11):6309–6313, 1980.
- esi et al. an algorithm for finding conserved secondary structure motifs in unaligned sequences. *NAR*, 32:3258–3269, 2004.
- et al. *PNAS*, 102:4006–4009, 2005.
- . mRNA stability in mammalian cells. *Microbiol Rev*, 59(3):423–450, 1995.
- ot. Spelling approximate or repeated motifs using a suffix tree. *LNCS*, 111:11–127, 1998.
- th et al. *Eur. J. Pharm. Sci.*, 11:191, 2000.
- g et al. a framework for RNA silencing in plants. *Genes Dev*, 17:49–63, 2003.
- co et al. *Nature New Biology*, 246:40–41, 1973.
- nderzande. *Lattice Models of Polymers (Cambridge Lecture Notes in Physics 11)*. Cambridge University Press, 1998.
- terman and T. Smith. Rapid dynamic programming algorithms for RNA secondary structure. *Adv. Appl. Math.*, 7:455–464, 1986.
- sh, N. Scherberg, R. Gilmore, and D. Steiner. Translational control of insulin biosynthesis. *Biochem. J.*, 235:459–467, 1986.
- kie, K. Dickson, and N. Gray. Regulation of mRNA translation by 5'- and 3'-untranslated factors. *Trends Biochem Sci*, 28:182–188, 2003.
- g et al. Decay rates of human mRNAs: correlation with functional characteristics and genomic attributes. *Genome Res*, 13:1863–1872, 2003.
- erstein, E. Eskin, and Z. Yakhini. Sequence motifs in ranked expression data. In *The RECOMB Satellite Workshop on Regulatory Genomics*, 2004.
- erstein, M. Ziv-Ukelson, R. Y. Pinter, and Z. Yakhini. A high-throughput approach for associating microRNAs with their activity conditions. In *RECOMB*, 133–151, 2005.
- giaga, J. Belasco, and M. Greenberg. The nonamer uuauuuuu is the key au-rich sequence motif that mediates mRNA degradation. *Mol. Cell. Biol.*, 15:2219–2230, 1995.
- ker. Computer prediction of RNA structure. *Methods Enzymol.*, 180:262–288, 1989.
- ker. *NAR*, (13):3406–15, 2003.
- ker and P. Stiegler. Optimal computer folding of large RNA sequences using thermodynamic and auxiliary information. *NAR*, 9(1):133–148, 1981.