

Sparsification of RNA Structure Prediction Including Pseudoknots

Mathias Möhl^{1*}, Raheleh Salari^{2*}, Sebastian Will^{1,3*},
Rolf Backofen^{1**} and S. Cenk Sahinalp^{2**}

¹ Bioinformatics, Institute of Computer Science, Albert-Ludwigs-Universität,
Freiburg, Germany

² Lab for Computational Biology, School of Computing Science, Simon Fraser
University, Burnaby, BC, Canada

³ Computation and Biology Lab, CSAIL, MIT, Cambridge MA, USA

Abstract. Although many RNA molecules contain pseudoknots, computational prediction of pseudoknotted RNA structure is still in its infancy due to high running time and space consumption implied by the dynamic programming formulations of the problem. In this paper, we introduce sparsification to significantly speedup the dynamic programming approaches for pseudoknotted RNA structure prediction, which also lower the space requirements. Although sparsification has been applied to a number of RNA-related structure prediction problems in the past few years, we provide the first application of sparsification to pseudoknotted RNA structure prediction specifically and to handling gapped fragments more generally - which has a much more complex recursive structure than other problems to which sparsification has been applied. We show that sparsification, when applied to the fastest, as well as the most general pseudoknotted structure prediction methods available, - respectively the Reeder-Giegerich algorithm and the Rivas-Eddy algorithm - reduces the number of "candidate" substructures to be considered significantly. In fact, experimental results on the sparsified Reeder-Giegerich algorithm suggest a linear speedup over the unparsified implementation.

1 Introduction

Recently discovered catalytic and regulatory RNAs [1, 2], exhibit their functionality due to specific secondary and tertiary structures [3, 4]. The vast majority of computational analysis of non-coding RNAs have been restricted to nested secondary structures, neglecting pseudoknots - which are "among the most prevalent RNA structures" [5]. For example, Xayaphoummine et al. [6] estimated that up to 30% of the base pairs in G+C-rich sequences form pseudoknots.

However the general problem of pseudoknotted RNA structure prediction is NP-hard. As a result, a number of approaches have been introduced for handling restricted classes of pseudoknots [7–13]. Condon *et al.* [14] give an overview of

* joint first authors

** to whom correspondence should be addressed

their structure classes and the algorithm-specific restrictions and Möhl *et al.* [15] develop a general framework showing that all these algorithms follow a general scheme, which they use for efficient alignment of pseudoknotted RNA.

The most general algorithm (with respect to the pseudoknot classes handled) among the above by Rivas and Eddy (R&E) has a running time of $O(n^6)$ time and space consumption of $O(n^4)$. It is therefore too expensive to directly apply this algorithm for large scale data analysis. Unfortunately, even the most efficient algorithm by Reeder and Giegerich (R&G) still has a high running time of $O(n^4)$, although it strongly restricts the class of predictable pseudoknots.

In this paper we introduce the technique of sparsification to the problem of pseudoknotted RNA structure prediction. Sparsification improves the expected running time and space usage of a dynamic programming based structure prediction algorithm without introducing additional restrictions on the structure class handled or compromising the optimality of solutions. Sparsification has been recently applied to improve time and space complexity of various existing RNA-related structure prediction algorithms. In particular, it turned out to be successful for RNA folding for pseudoknot-free structures [16, 17], simultaneous alignment and folding [18] as well as RNA RNA interaction prediction [19].

Contributions. We study sparsification of pseudoknotted RNA structure prediction. Algorithms developed for this problem differ from the previously sparsified algorithms by their use of gapped fragments and their more complex recursion structure. Our main contribution in this paper is the solution to the algorithmic challenges due to this increased complexity. Among all DP based pseudoknot prediction algorithms, we focus on the fastest algorithm (R&G) and the most general one (R&E) and develop sparse variants of these dynamic programming algorithms. Due to sparsification, the resulting algorithms need to consider only a limited number of candidates substructures compared to the original algorithms. As a result, we analyze the theoretical worst case complexities in terms of the number of candidate substructures. We also present experimental results, comparing our implementations of the original and sparsified *R&G* algorithm. These results suggest a significant (roughly a linear factor) reduction in the number of candidates over the original algorithm.

2 Sparsification of the Reeder and Giegerich algorithm

The R&G algorithm [13] predicts the minimum free energy structure allowing canonical pseudoknots for a sequence S of length n . It extends the Zuker algorithm by adding one more matrix K (for knot), where $K(i, j)$ denotes the energy for the best *canonical* pseudoknot that starts at position i and ends at position j .⁴ Canonical pseudoknots are defined as follows. Each pair of base pairs $p_1 = (i, i')$ and $p_2 = (j', j)$ with $i < j' < i' < j$ induces one canonical pseudoknot that consists of two crossing stems $\{(i, i'), (i + 1, i' - 1), \dots, (i + d_{i,i'} - 1, i' - d_{i,i'} + 1)\}$

⁴ The original presentation of the algorithm in terms of the ADP framework does not explicitly consider a matrix K but only a motif *knot*

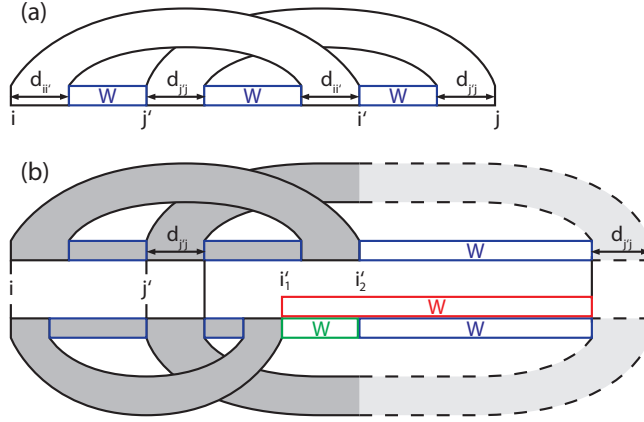


Fig. 1. Recursion for canonical pseudoknots (a) and their sparsification (b).

and $\{(j', j), (j' + 1, j - 1), \dots, (j' + d_{j',j} - 1, j - d_{j',j} + 1)\}$ where the stacking length of the two stems, $d_{i,i'}$ and $d_{j',j}$, respectively, is chosen as large as possible such that still all base pairs are valid Watson-Crick base pairs.

To allow for sparsification, we restrict the scoring scheme slightly such that the energy of a canonical pseudoknot only depends on the left ends of its base pairs⁵ and hence can be described as PK-Energy $(i, d_{i,i'}, j', d_{j',j})$. Then,

$$K(i, j) = \min_{i', j'} \text{score}(i, j', i', j) \quad (1)$$

with $\text{score}(i, j', i', j) =$

$$\left(\begin{array}{l} \text{PK-Energy}(i, d_{i,i'}, j', d_{j',j}) + \\ W(i + d_{i,i'}, j' - 1) + W(j' + d_{j',j}, i' - d_{i,i'}) + W(i' + 1, j - d_{j',j}) \end{array} \right). \quad (2)$$

As shown in Fig. 1(a), for each canonical pseudoknot starting at i and ending at j the recursion decomposes into the pseudoknot itself and the three fragments in-between its two crossing stems. Such pseudoknots add one case in the computation of a matrix entry $W(i, j)$, which, as in the Zuker algorithm, contains the optimal energy of a substructure starting at position i and ending at position j . Due to the restriction to canonical pseudoknots, the recursion of R&G minimizes only over all possible instances of i' and j' , because the maximal stacking lengths $d_{i,i'}$ and $d_{j',j}$ are uniquely determined once i' and j' are fixed. Furthermore, Reeder and Giegerich note that the maximal stacking length $d_{x,y}$ can be precomputed for all x, y in $O(n^3)$ time and stored in an $O(n^2)$ table.

In order to sparsify the algorithm, we develop an appropriate notion of a *candidate* such that it is not necessary to minimize over all possible i' and j' but only over the candidates.

⁵ The restricted scoring scheme does not distinguish between G-C and G-U base pairs in pseudoknot-stems, since their left ends are identical.

Definition 1 (R&G candidate).

Let $i < j' < i'_1 < i'_2$ and $d_{j',j} \leq i'_1 - j'$. Then i'_1 dominates i'_2 with respect to $(i, j', d_{j',j})$ iff

$$\text{score}_{i'_2}(i, j', i'_2) \geq \text{score}_{i'_1}(i, j', i'_1), \text{ where}$$

$$\begin{aligned} \text{score}_{i_c}(i, j', i') &:= \text{PK-Energy}(i, d_{i,i'}, j', d_{j',j}) \\ &+ W(i + d_{i,i'}, j' - 1) + W(j' + d_{j',j}, i' - d_{i,i'}) + W(i' + 1, i'_c). \end{aligned}$$

We say that i'_2 is a candidate with respect to $(i, j', d_{j',j})$ if there does not exist any i'_1 that dominates it.

The notion of a candidate is visualized in Fig. 1(b). There, i'_1 dominates i'_2 if the score for the gray area at the top (including the dashed part whose exact position is not determined) is not better than the score for the corresponding gray area at the bottom plus the green part. Note that these scores (and hence the candidate i') depend only on i, j' , and $d_{j',j}$ and are independent of $d_{i,i'}$ and j . The following lemma shows that the notion of a candidate given in Def. 1 is suitable for sparsification, i.e. some i' needs to be considered in the recursion (for all j) only if it is a candidate, because otherwise it is dominated by a candidate that yields a better score.

Lemma 1 (R&G sparsification). Let i'_2 be dominated by i'_1 with respect to some $(i, j', d_{j',j})$. Then for all j it holds $\text{score}(i, j', i'_1, j) \leq \text{score}(i, j', i'_2, j)$.

Proof. We start with the inequality of Def. 1 and add $W(i'_2 + 1, j - d_{j',j})$ on both sides. Then the claim follows immediately from $W(i'_1 + 1, j - d_{j',j}) \leq W(i'_1 + 1, i'_2) + W(i'_2 + 1, j - d_{j',j})$. In Fig. 1(b) this corresponds to the fact that the score for the red box is at least as good as the score from the green and the blue box together. This triangle inequality holds by the correctness of the (unsparsified) algorithm: For all $x < y < z$ we have $W(x, y) + W(y + 1, z) \leq W(x, z)$ since the concatenation of the best structures for the ranges (x, y) and (y, z) always forms a valid structure for the range (x, z) with score $W(x, y) + W(y + 1, z)$ which is hence never better than the optimal score $W(x, z)$ for that range. \square

The sparsified algorithm maintains lists L_i of candidates for each pair $(j', d_{j',j})$ since only the lists for one i need to be maintained in memory at the same time. Whenever in the computation of some $\text{score}(i, j', i', j)$ the i' is considered the first time for this i and j' , it is checked whether it is a candidate and if so, it is added to the respective list. For all other instances of j, i' is then considered only if it is contained in the list. The sparsified algorithm is given by the following pseudo-code ($n := |S|$).

```

1: for  $i := n$  to 1 do
2:   for all  $d_{j',j}, j' \leq n$  do  $L_i(j', d_{j',j}) :=$ empty list;
3:   for  $j := i + 3$  to  $n$  do
4:      $K(i, j) := \infty$ 
5:     for  $j' := i + 1$  to  $j - 2$  do
6:       // check new elements for candidacy
```

```

7:   for  $i_c := \max\{j' + d_{j',j}, \text{checked}_{i,j',d_{j',j}} + 1\}$  to  $j - d_{j',j}$  do
8:     if  $\text{score}_{i_c}(i, j', i_c) < \text{score}_{i_c}(i, j', i')$  for all  $i' \in L_i(j', d_{j',j})$  then
9:       add  $i_c$  to  $L_i(j', d_{j',j})$ 
10:    end if
11:  end for
12:   $\text{checked}_{i,j',d_{j',j}} := \max(\text{checked}_{i,j',d_{j',j}}, j - d_{j',j})$ 
13:  // iterate over all candidates
14:   $K_{i,j',j} := \infty$ 
15:  for all  $i' \in L_i(j', d_{j',j})$  do
16:     $K_{i,j',j} := \min\{K_{i,j',j}, \text{score}(i, j', i', j)\}$ 
17:  end for
18:   $K(i, j) := \min\{K(i, j), K_{i,j',j}\}$ 
19: end for
20: compute matrix entries  $V(i, j)$  and  $W(i, j)$  as in Wexler et al.
21:  $W(i, j) := \min(W(i, j), K(i, j))$ 
22: end for
23: end for

```

The candidate lists are initialized in line 2. In lines 7 to 11 all new values i_c that have not been considered so far, are tested for candidacy. Here, $\text{checked}_{i,j',d_{j',j}}$ denotes the largest i' that has been checked for candidacy in list $L_i(j', d_{j',j})$.

Lines 14 to 17 compute scores $\text{score}(i, j', i', j)$ for all candidates i' . In line 20, we compute $W(i, j)$ and $V(i, j)$ as in the sparsified pseudoknot-free structure prediction approach due to Wexler *et al.* [16]. The computation of matrices K and W is interleaved such that all entries $K(i, j)$ and $W(i, j)$ are computed before all entries $K(i', j')$ and $W(i', j')$ for $i \leq i' \leq j' \leq j$ and $i \neq i'$ or $j \neq j'$.

Complexity Analysis Whereas the original algorithm requires $O(n^4)$ time (for $n = |S|$), the sparsified variant requires $O(n^3L)$ time where L is the total size for all candidate lists of some i i.e. $L := \max_i \sum_{j', d_{j',j}} |L_i(j', d_{j',j})|$. Obviously, $L \leq n$. In order to maintain the asymptotic space complexity $O(n^2)$ of the original algorithm, we do not maintain all lists $L_i(j', d_{j',j})$ in memory but only the lists with $d_{j',j} \leq k$ where $k > 0$ is a small constant. Please note that to keep presentation simple, we didn't make this explicit in the pseudo-code. Since the maximal stacking length is usually small, there are only very few instances of j with $d_{j',j} > k$ such that for those few j it is cheap to consider all i' as candidates. Hence, we store $O(kn) = O(n)$ candidate lists each requiring at most $O(n)$ space.

3 Sparsification of the Rivas and Eddy Algorithm

The class of structures predicted by the R&E algorithm [8], here called class of R&E structures, is the most general RNA secondary structure prediction algorithm described in the literature [14]. To keep presentation simple we explain the sparsification strategy for a base-pair maximization algorithm that handles the R&E structure class. Finally, we motivate that sparsification can be transferred to the R&E energy minimization algorithm.

First, we give recursions of base pair maximization for R&E structures. Note that the recursions are intentionally very close to the recursions of the R&E energy minimization algorithm. After initialization for $i \geq j$ and $k \geq l$

$$W(i, j) = \begin{cases} 0 & \text{if } i = j \text{ or } i = j + 1 \\ -\infty & \text{if } i > j + 1 \end{cases} \quad \text{and} \quad \begin{cases} W(i, j; k, l) = -\infty & \text{if } j < i \text{ or } l < k \\ W(i, i; k, k) = \text{bp}(i, k) \end{cases}$$

where $\text{bp}(i, j) = \begin{cases} 1 & \text{if } S_i, S_k \text{ complementary} \\ -\infty & \text{otherwise,} \end{cases}$ is the *base pair contribution*,

the recursions (R&E recursions) are given for $1 \leq i < j < k < l \leq |S|$ as

$$W(i, j) = \max \begin{cases} W(i, j - 1) & (12') \\ \text{bp}(i, j) + W(i + 1, j - 1) & (1'21') \\ \max_{j'} W(i, j' - 1) + W(j', j) & (12) \\ \max_{j', k', l'} W(i, j' - 1; k' + 1, l' - 1) + W(j', k'; l', j) & (1212) \end{cases}$$

$$W(i, j; k, l) = \max \begin{cases} W(i + 1, j; k, l) & (1'2G2) \\ W(i, j - 1; k, l) & (12'G1) \\ W(i, j; k + 1, l) & (1G2'1) \\ W(i, j; k, l - 1) & (1G12') \\ \max_{j'} W(i, j') + W(j' + 1, j; k, l) & (12G2) \\ \max_{j'} W(i, j' - 1, j; k, l) + W(j', j) & (12G1) \\ \max_{l'} W(i, j; l' + 1, l) + W(k, l') & (1G21) \\ \max_{l'} W(i, j; k, l' - 1) + W(l', l) & (1G12) \\ \max_{j', k'} W(i, j' - 1; k' + 1, l) + W(j', j; k, k') & (12G21) \\ \max_{j', k'} W(i, j' - 1; k, k' - 1) + W(j', j; k', l) & (12G12) \\ \max_{k', l'} W(i, j; k' + 1, l' - 1) + W(k, k'; l', l) & (1G212) \\ \max_{i', j'} W(i, i' - 1; j' + 1, j) + W(i', j'; k, l) & (121G2). \end{cases}$$

It is easy to check that $W(1, |S|)$ is the maximal number of base pairs in a R&E structure of S , because the recursions perform the same decompositions as the original R&E recursions. Note that $W(i, j; k, l)$ is the maximal number of base pairs in structures with at least one base pair that spans the gap. We label each recursion case in a way that illustrates the type of the decomposition of this case. The idea of these labels is taken from Möhl *et al.* [15], where we developed a type system for decompositions, which there are called splits. For this reason, we call these labels split types, however, we won't need any details of the typing system. The decomposition by R&E is illustrated in Figure 2.

A *fragment* is defined as a set of positions of the fixed sequence S . The fragments corresponding to matrix entries in the R&E recursion can be described conveniently by their boundaries. We distinguish *ungapped fragments* $F = \{i, \dots, j\}$, written (i, j) , and *1-gap fragments* $F' = \{i, \dots, j\} \cup \{k, \dots, l\}$, written $(i, j; k, l)$ where i, j, k, l , are called *boundaries* of respective F or F' . A *split* of a fragment F is a tuple (F_1, F_2) such that $F = F_1 \cup F_2$ and $F_1 \cap F_2 = \emptyset$.

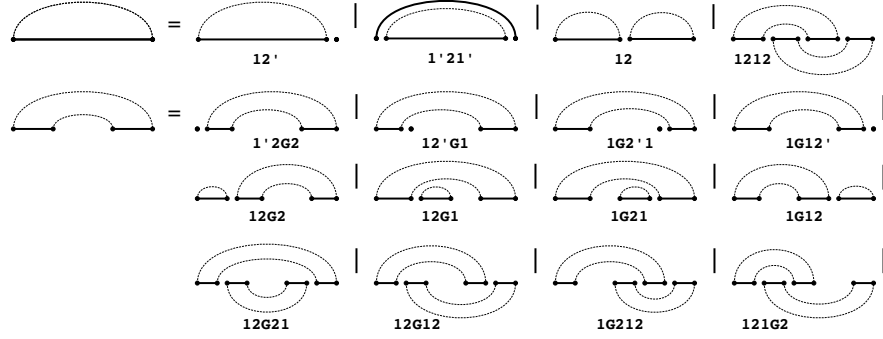


Fig. 2. Decomposition for R&E base pair maximization annotated with labels, i.e. split types, of the corresponding recursion cases.

For our sparsification approach, we will show that in each recursion case, certain optimally decomposable fragments do not have to be considered for computing an optimal solution, because each decomposition using these fragments can be replaced by a decomposition using a smaller fragment. We define optimal decomposability with respect to the split type of a R&E recursion case.

Definition 2 (Optimally decomposable). A fragment F is optimally decomposable by a split of type T (T -OD) iff there is a split (F_1, F_2) that occurs in recursion case T and $W(F_1) + W(F_2) \geq W(F)$.

A fragment F is optimally decomposable w.r.t a set of split types \mathcal{T} (\mathcal{T} -OD) iff F is T -OD for some $T \in \mathcal{T}$.

Here, we emphasize that testing T -OD for a fragment F is simple in a run of the DP algorithm. After evaluating the case T in the computation of $W(F)$, one compares the maximum of the case to $W(F)$. For example, a fragment $(i, j; k, l)$ is 12G21-OD iff $W(i, j; k, l) = \max_{j', k'} W(i, j' - 1; k' + 1, l) + W(j', j; k, k')$.

In the following we show that for the maximization in a recursion case T , we do not need to consider T' -OD fragments as second fragment of the split, where T' is from a T -specific set of split types. As an example consider the recursion case 12G21, which splits fragments $(i, j; k, l)$ into $F_1 = (i, j' - 1; k' + 1, l)$ and $F_2 = (j', j; k, k')$. Assume that F_2 is 12G21-OD. Then we can show that every evaluation of $W(F)$ where $W(F) = W(F_1) + W(F_2)$ can be replaced by another at least equally good evaluation that splits F into F'_1 and $F'_2 \subset F_2$, where F'_2 is the second fragment in the 12G21-split of F_2 . However, note that the argument is split type specific and cannot be applied e.g. when F_2 is 12G12-OD.

For sparsifying R&E, we define the following sets of split types.

$$\begin{aligned}
\mathcal{T}_{12}^{\text{RE}} &= \{12\} & \mathcal{T}_{1212}^{\text{RE}} &= \{12G2, 12G1, 1G21\} \\
\mathcal{T}_{12G1}^{\text{RE}} &= \mathcal{T}_{1G12}^{\text{RE}} = \mathcal{T}_{1G21}^{\text{RE}} = \{12\} & \mathcal{T}_{12G2}^{\text{RE}} &= \{12G2\} \\
\mathcal{T}_{12G21}^{\text{RE}} &= \{12G2, 1G12, 12G21\} & \mathcal{T}_{12G12}^{\text{RE}} &= \{12G2, 1G21, 12G12\} \\
\mathcal{T}_{1G212}^{\text{RE}} &= \{12G1, 1G21, 12G21\} & \mathcal{T}_{121G2}^{\text{RE}} &= \{12G2, 12G1, 121G2\}
\end{aligned}$$

These sets are defined such that in a recursion case T , whenever the second fragment of a split (F_1, F_2) of F can be optimally decomposed by a split of a type in $\mathcal{T}_T^{\text{RE}}$, a different split (F'_1, F'_2) of type T can be applied to F , where $F'_2 \subset F_2$. As we show later, this split will be just as good as (F_1, F_2) for computing $W(F)$.

Then, one systematically obtains sparsified recursion equations $W'(i, j)$ and $W'(i, j; k, l)$ from the equations for $W(i, j)$ and $W(i, j; k, l)$ by replacing symbol W by W' and modifying them in the following way. For each case T in the recursion of $W(i, j)$ and $W(i, j; k, l)$ that maximizes over $W(F_1) + W(F_2)$ for respective splits of the fragment $F = (i, j)$ or $F = (i, j; k, l)$, maximize only over fragments F_2 that are not $\mathcal{T}_T^{\text{RE}}$ -OD. In an algorithm that evaluates the sparsified recursion, such non- $\mathcal{T}_T^{\text{RE}}$ -OD fragments correspond to entries of candidate lists. For example, case 12G21 of W is modified in the equation for $W'(i, j; k, l)$ to

$$\max_{j', k', (j', j; k, k') \text{ not } \mathcal{T}_{12G21}^{\text{RE}}\text{-OD}} W'(i, j' - 1; k' + 1, l) + W'(j', j; k, k') \quad (12G21 \text{ of } \mathbf{W}').$$

Theorem 1. *Let W be the matrix of the R&E recursion and W' its sparsified variant, then $W(1, |S|) = W'(1, |S|)$.*

Proof. We show for all $1 \leq i, j, k, l \leq |S|$, $W(i, j) = W'(i, j)$ and $W(i, j; k, l) = W'(i, j; k, l)$. First note that it holds that $W(i, j) \geq W'(i, j)$ and $W(i, j; k, l) \geq W'(i, j; k, l)$. The claim is shown by induction on the fragment size and a case distinction over recursion cases. For the case of split type 12, we show that

$$\max_{j'} W(i, j' - 1) + W(j', j) = \max_{j', (j', j) \text{ not } \mathcal{T}_{12}^{\text{RE}}\text{-OD}} W'(i, j' - 1) + W'(j', j).$$

Let (j', j) be 12-OD for some $j' : i \leq j' \leq j$. By IH, it suffices to find a (smaller) fragment (j'', j) , where $j'' > j'$ and $W(i, j'' - 1) + W(j'', j) \geq W(i, j' - 1) + W(j', j)$. Either (j', j) is not 12-OD or there is a j'' , such that $W(j', j) = W(j', j'' - 1) + W(j'', j)$ and thus $W(i, j'' - 1) + W(j'', j) \geq W(i, j' - 1) + W(j', j)$ because

$$\begin{aligned} W(i, j'' - 1) + W(j'', j) &\geq_{\Delta\text{-ineq}} W(i, j' - 1) + W(j', j'' - 1) + W(j'', j) \\ &=_{12\text{-OD}} W(i, j' - 1) + W(j', j). \end{aligned}$$

The triangle inequality (Δ -ineq) is an immediate consequence of the correctness of the recursion for W . Thus, for the decompositions of all recursion cases there holds such a corresponding inequation. Analogous arguments can be given for all other modified recursion cases. Exemplarily, we elaborate the argument for the complex case 12G21. Let $F_1 = (i, j' - 1; k' + 1, l)$ and $F_2 = (j', j; k, k')$, such that (F_1, F_2) is a split of type 12G21 of $(j, j; k, k')$. We need to show for all $\mathcal{T}_{12G21}^{\text{RE}}$ -OD fragments F_2 there are non-empty ungapped or 1-gap fragments F'_1 and F'_2 , where $F'_1 \cup F'_2 = F_2$, $F'_1 \cap F'_2 = \emptyset$, and $W(F_1 \cup F'_1) + W(F'_2) \geq W(F_1) + W(F_2)$ and the split $(F_1 \cup F'_1, F'_2)$ occurs in a recursion case of R&E. Again, either F_2 is not $\mathcal{T}_{12G21}^{\text{RE}}$ -OD or one of the following cases applies. Case 1 (12G2): for some j'' , $W(j', j; k, k') = W(j', j'' - 1) + W(j'', j; k, k')$. Then, the claim holds for $F'_1 = (j', j'' - 1)$ and $F'_2 = (j'', j; k, k')$ by triangle inequality and split $(F_1 \cup F'_1, F'_2)$ occurs in recursion case 12G21. Case 2 (2G21): for some k'' , $W(j', j; k, k') =$

$W(j', j; k, k'') + W(k'' + 1, k')$. The claim holds for $F'_2 = (j', j; k, k'')$. Case 3 (12G21): for some j'', k'' , $W(j', j; k, k') = W(j', j''-1; k''+1, k') + W(j'', j; k, k')$. Again, this satisfies the claim by triangle inequality. \square

Algorithm The recursion equation W' tailors a sparsified dynamic programming algorithm for the evaluation of $W'(1, |S|)$ with very limited overhead. We maintain separate candidate lists for each sparsified recursion case. As already mentioned, the T -OD properties of each fragment F can be easily checked after evaluation of each case of $W(F)$. A fragment is added to a candidate list for recursion case T iff it is not $\mathcal{T}_T^{\text{RE}}$ -OD. The maximizations are restricted to run only over the candidates in the respective candidate list. Their intended use dictates the exact nature of such candidate lists. For a case T , which splits a fragments T into T_1 and T_2 , there are candidate lists for all boundaries of a fragment T_2 that are not adjacent to boundaries of T_1 due to split type T . The list entries are tuples of the adjacent boundaries and the fragment score for T_2 . In order to profit from a reduced number of candidates in space, we maintain two three-dimensional slices of the matrix for $W(i, j; k, l)$, storing entries only for the current i and $i + 1$. Scores $W(i, j; k, l)$ for larger i are stored for candidates only.

R&E Free Energy Minimization Sparsification is analogously applied to the energy minimizing R&E algorithm. This algorithm distinguishes several additional matrices that contain minimal energies for fragments (i, j) or $(i, j; k, l)$ under the condition that respectively the base pair (i, j) or base pairs (i, l) and (j, k) or one of them exist. Almost all decompositions in the recursion for these matrices are of discussed split types and are sparsified analogously. The only notable exception is due to internal loops. Internal loops require minimizing over all possible positions of the inner loop base pair, where commonly the loop size is restricted by a constant K such that minimizing takes constant time. However, handling inner loops requires access to entries of non-candidate fragments $(i', j'; k', l')$ for $i \leq i' \leq i + K + 2$. This is handled by maintaining matrix slices for i to $i + K + 2$ in $O(n^3)$ space, which preserves total space complexity.

Complexity Analysis The described algorithm profits from sparsification in time and space. Compared to $O(n^6)$ time and $O(n^4)$ space of the unsparsified algorithm (for $n = |S|$), we obtain complexities in the number of candidates. Let Z_T denote the maximal length of a candidate lists for case T and Z denote the total number of entries in all lists. Then, the time complexity is $O(n^2(Z_{12} + Z_{1212}) + n^4(Z_{12G2} + Z_{12G1} + Z_{1G21} + Z_{1G12} + Z_{12G21} + Z_{12G12} + Z_{1G212} + Z_{121G2}))$ and space complexity is $O(n^3 + Z)$. In the worst case, Z_{12} , Z_{12G2} , Z_{12G1} , Z_{1G21} and Z_{1G12} are $O(n)$, Z_{12G21} , Z_{12G12} , Z_{1G212} , Z_{121G2} are $O(n^2)$, and Z_{1212} is $O(n^3)$; finally Z is $O(n^4)$ in the worst case.

4 Experimental Results

In order to evaluate the effect of sparsification on pseudoknotted RNA secondary structure prediction, we implemented original and sparsified variants of the Reeder and Giegerich (R&G) algorithm.

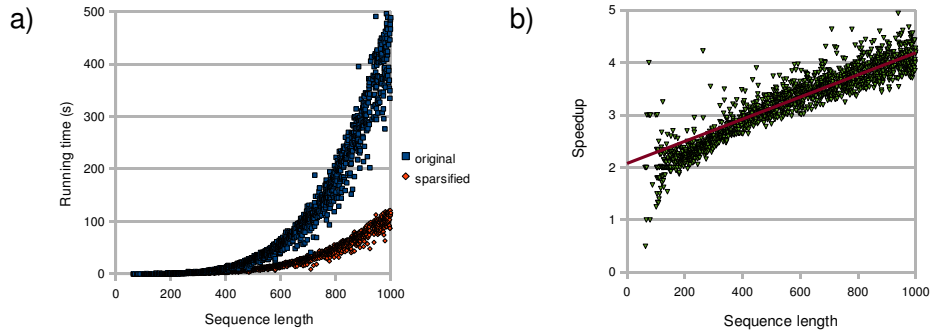


Fig. 3. Running times of the original and sparsified variants of the R&G algorithm.

Data Set We obtained all RNA sequences from PseudoBase[20], which are known to have some pseudoknots in their secondary structures. This set contains 294 sequences that their length is distributed between 76nt and 93399nt. We randomly divided all long sequences into subsequences shorter than 1000nt. Therefore the data set that we used in our experiments contains 1563 sequences with length between 76nt and 1000nt.

Performance We applied both variants of the R&G algorithm to our data set. Fig. 3 shows the running time of the algorithms on a server with Intel Core Duo CPU at 2.53GHz and 4GB RAM. The results in Fig. 3 show that sparsification significantly improves the running time of the R&G algorithm. As the RNA sequences get longer, the relative performance of the sparsified algorithm (with respect to the non-sparsified ones) improves. Fig. 3.(b) shows the speedup of the sparsified algorithm, which fits well to a linear regression ($R^2 = 0.84$).

Number of candidates For a better understanding of the effect of sparsification on the R&G algorithm, we measured the number of (i', j') pairs which are checked in each fragment $[i, j]$ in both original and sparsified variants of the algorithm. Note that the number of (i', j') pairs is in order of $O((j - i)^2)$ in the worst case. Fig. 4 shows the average number of (i', j') pairs on fragments of equal length which are checked by the two variants of the algorithm. As expected, this amount is significantly smaller for the sparsified algorithm compared to the original one. Moreover, we observe that as the fragments get longer, the difference between the average number of (i', j') pairs in the sparsified and the original algorithm increases. We define the work load per each fragment $[i, j]$ as the number of candidate (i', j') pairs. Figure 4(b), shows a significant reduction of the work load in the sparsified algorithms. As it can be seen for subsequences of length 1000nt, the work load by the sparsified algorithm is reduced by a factor of about 10 compared to the original algorithm. Note that the work load reduction at fragment length 1000nt does not yield the same speedup for sequences of length 1000nt (here this speedup is about 3.5, confer Fig.3(b)), because for a sequence of length n , all fragments of smaller length are processed by the algorithm.

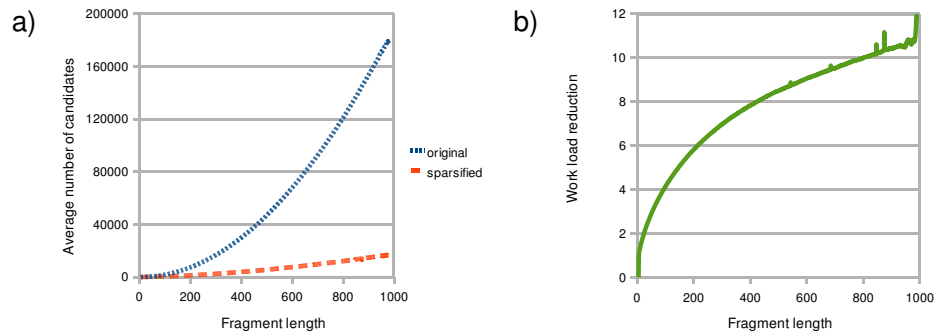


Fig. 4. Average number of (i', j') candidates in the original and sparsified variants of the R&G algorithm.

5 Conclusion

The presented work gives two examples for sparsification in the context of gap fragments and a complex recursion structure. Since we successfully sparsified the fastest and the most complex pseudoknot structure prediction algorithm for RNA, it is likely that all other DP-based pseudoknot-algorithm can be sparsified. Thus, the paper motivates further generalization of sparsification for systematic application to complex DP-algorithms as RNA structure prediction algorithms. Even more, by providing detailed examples the paper directly prepares such generalization. Our results from an implementation of the sparsified Reeder and Giegerich algorithm show a significant, presumably even linear, expected work load reduction due to sparsification.

Acknowledgments This work is partially supported by DFG grants WI 3628/1-1 and BA 2168/3-1 R. Salari was supported by SFU-CTEF funded Bioinformatics for Combating Infectious Diseases Project co-lead by Sahinalp. S.C. Sahinalp was supported by MITACS, NSERC, the CRC program and the Michael Smith Foundation for Health Research.

References

1. Sharp, P.A.: The centrality of RNA. *Cell* **136**(4) (2009) 577–80
2. Amaral, P.P., Dinger, M.E., Mercer, T.R., Mattick, J.S.: The eukaryotic genome as an RNA machine. *Science* **319**(5871) (2008) 1787–9
3. Washietl, S., Pedersen, J.S., Korbil, J.O., Stocsits, C., Gruber, A.R., Hackermuller, J., Hertel, J., Lindemeyer, M., Reiche, K., Tanzer, A., Ucla, C., Wyss, C., Antonarakis, S.E., Denoeud, F., Lagarde, J., Drenkow, J., Kapranov, P., Gingeras, T.R., Guigo, R., Snyder, M., Gerstein, M.B., Reymond, A., Hofacker, I.L., Stadler, P.F.: Structured RNAs in the ENCODE selected regions of the human genome. *Genome Res* **17**(6) (2007) 852–64
4. Mattick, J.S., Makunin, I.V.: Non-coding RNA. *Hum Mol Genet* **15 Spec No 1** (2006) R17–29

5. Staple, D.W., Butcher, S.E.: Pseudoknots: RNA structures with diverse functions. *PLoS Biol* **3**(6) (2005) e213
6. Xayaphoummine, A., Bucher, T., Thalmann, F., Isambert, H.: Prediction and statistics of pseudoknots in RNA structures using exactly clustered stochastic simulations. *Proc. Natl. Acad. Sci. USA* **100**(26) (2003) 15310–5
7. Lyngso, R.B., Pedersen, C.N.S.: Pseudoknots in RNA secondary structures. In: *Proc. of the Fourth Annual International Conferences on Computational Molecular Biology (RECOMB'00)*, ACM Press (2000) BRICS Report Series RS-00-1.
8. Rivas, E., Eddy, S.R.: A dynamic programming algorithm for RNA structure prediction including pseudoknots. *Journal of Molecular Biology* **285**(5) (1999) 2053–68
9. Uemura, Y., Hasegawa, A., Kobayashi, S., Yokomori, T.: Tree adjoining grammars for RNA structure prediction. *Theoretical Computer Science* **210** (1999) 277 – 303 Paper as Print Copy.
10. Akutsu, T.: Dynamic programming algorithms for RNA secondary structure prediction with pseudoknots. *Discrete Applied Mathematics* **104** (2000) 45–62
11. Deogun, J.S., Donis, R., Komina, O., Ma, F.: RNA secondary structure prediction with simple pseudoknots. In: *APBC '04: Proceedings of the second conference on Asia-Pacific bioinformatics*, Darlinghurst, Australia, Australia, Australian Computer Society, Inc. (2004) 239–246
12. Dirks, R.M., Pierce, N.A.: A partition function algorithm for nucleic acid secondary structure including pseudoknots. *J Comput Chem* **24**(13) (2003) 1664–77
13. Reeder, J., Giegerich, R.: Design, implementation and evaluation of a practical pseudoknot folding algorithm based on thermodynamics. *BMC Bioinformatics* **5** (2004) 104
14. Condon, A., Davy, B., Rastegari, B., Zhao, S., Tarrant, F.: Classifying RNA pseudoknotted structures. *Theoretical Computer Science* **320**(1) (2004) 35–50
15. Möhl, M., Will, S., Backofen, R.: Lifting prediction to alignment of RNA pseudoknots. *Journal of Computational Biology* (2010) Accepted.
16. Wexler, Y., Zilberstein, C.B.Z., Ziv-Ukelson, M.: A study of accessible motifs and rna folding complexity. In Apostolico, A., Guerra, C., Istrail, S., Pevzner, P.A., Waterman, M.S., eds.: *Proc. of the Tenth Annual International Conferences on Computational Molecular Biology (RECOMB'06)*. Volume 3909 of *Lecture Notes in Computer Science.*, Springer (2006) 473–487
17. Backofen, R., Tsur, D., Zakov, S., Ziv-Ukelson, M.: Sparse RNA folding: Time and space efficient algorithms. In Kucherov, G., Ukkonen, E., eds.: *Proc. 20th Symp. Combinatorial Pattern Matching*. Volume 5577 of *LNCS.*, Springer (2009) 249–262
18. Ziv-Ukelson, M., Gat-Viks, I., Wexler, Y., Shamir, R.: A faster algorithm for RNA co-folding. In Crandall, K.A., Lagergren, J., eds.: *WABI 2008*. Volume 5251 of *Lecture Notes in Computer Science.*, Springer (2008) 174–185
19. Salari, R., Möhl, M., Will, S., Sahinalp, S.C., Backofen, R.: Time and space efficient RNA-RNA interaction prediction via sparse folding. In: *Proc. of RECOMB 2010*. (2010) Accepted.
20. van Batenburg, F.H., Gulyaev, A.P., Pleij, C.W., Ng, J., Oliehoek, J.: Pseudobase: a database with RNA pseudoknots. *Nucleic Acids Research* **28**(1) (2000) 201–4