

# Efficient Alignment of RNAs With Pseudoknots Using Sequence Alignment Constraints

Byung-Jun Yoon

## Abstract

When aligning RNAs, it is important to consider both the secondary structure similarity and primary sequence similarity to find an accurate alignment. However, algorithms that can handle RNA secondary structures typically have high computational complexity that limits their utility. For this reason, there have been a number of attempts to find useful alignment constraints that can reduce the computations without sacrificing the alignment accuracy. In this paper, we propose a new method for finding effective alignment constraints for fast and accurate structural alignment of RNAs, including pseudoknots. In the proposed method, we use a profile-HMM to identify the “seed” regions that can be aligned with high confidence. We also estimate the position range of the aligned bases that are located outside the seed regions. The location of the seed regions and the estimated range of the alignment positions are then used to establish the sequence alignment constraints. We incorporated the proposed constraints into the profile-csHMM (profile context-sensitive HMM) based RNA structural alignment algorithm. Experiments indicate that the proposed method can make the alignment speed up to 11 times faster without degrading the accuracy of the RNA alignment.

## Index Terms

RNA structural alignment, profile-csHMM (profile context-sensitive HMM), pseudoknot, sequence alignment constraint.

## I. INTRODUCTION

Sequence alignment lies at the heart of various computational methods that are used for analyzing biological sequences, such as RNAs and proteins. Alignment algorithms have been extensively used for comparing sequences to identify homologues, predict their structures, and infer their biological functions. Many functional noncoding RNAs (ncRNAs) are known to conserve their base-paired secondary structure as well as their primary sequence [1]. For this reason, when aligning RNAs, it is important to consider both structure and sequence similarities in order to find an accurate alignment that is biologically meaningful. For a similar reason, it is expedient to employ scoring schemes that can reasonably combine contributions from the secondary structure similarity as well as primary sequence similarity when performing an RNA similarity search, as it can significantly reduce the number of false positive predictions [2].

Conservation of the secondary structure gives rise to complicated symbol correlations between the pairing bases in the RNA sequence. Therefore, in order to take structural similarity into account, RNA alignment and search algorithms need to handle these base correlations in a principled manner. Until now, a number of probabilistic models have been proposed for this purpose [2], [3], where *stochastic context-free grammars* (SCFGs) and their variants have been especially popular. A typical problem of these models and the relevant algorithms is the high computational complexity. For example, the *Cocke-Younger-Kasami* (CYK) algorithm used in the SCFG-based alignment and search has a complexity of  $O(L^3)$ , where

Manuscript received June 17, 2008.

Byung-Jun Yoon is with the Department of Electrical and Computer Engineering, Texas A&M University, College Station, TX 77843-3128, USA.

$L$  is the length of the RNA that is to be aligned. Algorithms for simultaneous folding and alignment of RNAs (typically referred as the *Sankoff algorithm*) have an even higher complexity, which require  $O(L^{3N})$  computations for aligning  $N$  RNAs of length  $L$  [4]. The aforementioned algorithms do not consider pseudoknots<sup>1</sup>, and there will be a steep increase in complexity if we begin to consider RNA pseudoknots.

The high computational cost of RNA alignment and search algorithms limits their utility in practical applications, especially when the RNA of interest is long. To cope with this problem, there have been extensive research efforts to develop heuristic methods that can make these algorithms faster without degrading the accuracy. For example, let us consider the simultaneous folding and alignment algorithm. Its computational complexity is already  $O(L^6)$  for aligning just two RNAs, making it practically unusable for a larger number of RNAs. Even for pairwise alignments, the algorithm becomes quickly infeasible as the RNAs get longer. Therefore, in order to utilize these algorithms in practical applications, it is essential that we first reduce their computations. For this reason, most of the pairwise RNA alignment algorithms adopt various tricks to minimize the alignment time [5], [6], [7], [8], [9], [10], [11]. Similarly, a number of methods have been proposed to make RNA similarity searches faster, where the prescreening approach is a good example [12], [13], [14], [15]. The prescreening approach uses a simple model, such as a *profile-HMM* (*profile hidden Markov model*), to identify the regions that have a reasonable amount of (sequence) similarity. Only these regions will be passed to a more complex model, such as a *covariance model* (*CM*; *profile-SCFG*) [3] or a *profile-csHMM* (*profile context-sensitive HMM*) [16], [17], for further inspection. In fact, these are just a few examples, and there also exist other approaches for making RNA alignment and RNA search algorithms faster [18], [19].

Recently, we proposed an efficient RNA structural alignment algorithm based on profile-csHMMs, which can also be used for aligning RNAs that contain pseudoknots [17]. This algorithm finds the optimal alignment between a structured reference RNA and an unstructured target RNA, by taking both structure and sequence similarities into account. It was demonstrated that the profile-csHMM algorithm can find accurate alignment of RNA pseudoknots [17]. In this paper, we propose a novel method for finding effective sequence alignment constraints that can improve the computational efficiency of the profile-csHMM structural alignment algorithm. The overall organization of the paper is as follows. In Sec. II, we describe the concept of constrained alignment and briefly review some of the existing methods for finding the alignment constraints. After the review, we propose a new method for estimating the alignment constraints in Sec. III. Finally, Sec. IV describes how these constraints can be incorporated into the profile-csHMM based structural alignment method. Experimental results will be presented at the end of Sec. IV, which demonstrate the effectiveness of the proposed approach.

## II. CONSTRAINED SEQUENCE ALIGNMENT

Let us assume that we want to find the alignment of two RNAs  $\mathbf{X} = x_1x_2 \dots x_{L_1}$  (RNA-1) and  $\mathbf{Y} = y_1y_2 \dots y_{L_2}$  (RNA-2). The predicted sequence alignment can be uniquely represented by the set of aligned bases  $(x_m, y_n)$  in a matrix. For example, let us consider the RNA alignment in Fig. 1(a). The matrix shows

<sup>1</sup>RNA secondary structures with crossing base-pairs are called pseudoknots. Handling pseudoknots is computationally expensive, hence typically ignored by many algorithms.



As illustrated in the previous example, using appropriate sequence alignment constraints can greatly enhance the efficiency of the alignment algorithm. So, the natural question is how we can predict good alignment constraints that can minimize the alignment time without degrading the alignment accuracy. Until now, various methods have been proposed for restricting the alignment space to improve the efficiency of diverse RNA alignment and search algorithms [5], [6], [7], [8], [9], [11], [17], [18]. For example, the query-dependent banding (QDB) method [18] is used to make CM-based RNA alignment algorithms faster, by excluding the regions in the dynamic programming matrix that have insignificant probability. These regions can be precomputed based on the given CM and do not depend on the target database. *Foldalign* [7], an algorithm for simultaneous RNA folding and alignment, limits the maximum length of the RNA-motif as well as the maximum length difference between the subsequences that are being compared. Recent implementation of *Foldalign* [8] adopts a heuristic that prunes the dynamic programming matrix in order to reduce the overall time and memory requirements. Another RNA alignment and structure prediction algorithm called *Stemloc* [9] constrains the solution space by using “fold envelopes” and “alignment envelopes”. The fold envelopes are used to restrict the search over secondary structures and the alignment envelopes are used to restrict the possible alignments between the given sequences.

A recent implementation of *Dynalign* [11], a joint alignment and secondary structure prediction algorithm for two RNAs, assumes that the aligned bases in the respective RNAs should be located within a certain distance. To be more precise, the  $m$ -th base  $x_m$  in RNA-1 ( $\mathbf{X} = x_1x_2 \dots x_{L_1}$ ) can be aligned to the  $n$ -th base  $y_n$  in RNA-2 ( $\mathbf{Y} = y_1y_2 \dots y_{L_2}$ ), only if the following condition is satisfied

$$\left| \frac{L_2}{L_1}m - n \right| \leq M, \quad (1)$$

for a given  $M$ . For convenience, we refer this constraint as the  $M$ -constraint. The parameter  $M$  is used to specify the maximum distance between the alignable bases. By imposing this constraint, we are restricting the number of insertions and deletions in homologous sequences, which is a reasonable assumption for real biological sequences. The constrained alignment space is band-shaped as depicted in Fig. 2(a). Despite its simplicity, it has been shown that the proposed constraint works reasonably well [6], [11]. The latest implementation of *Dynalign* [6] takes a more principled approach for estimating the alignment region. In the new approach, a hidden Markov model (HMM) is used to estimate the set  $\mathcal{Q}$  of aligned positions  $(m, n)$ , whose *co-occurrence probability*<sup>2</sup>  $P(x_m \leftrightarrow y_n | \mathbf{X}, \mathbf{Y})$  is larger than a reasonably low threshold  $p_t$  ( $\approx 0$ )

$$\mathcal{Q} = \left\{ (m, n) \mid P(x_m \leftrightarrow y_n | \mathbf{X}, \mathbf{Y}) > p_t \right\}. \quad (2)$$

The estimated set  $\mathcal{Q}$  is used to constrain the final alignment. It was demonstrated that this technique can provide significant savings in computational time as well as a small improvement in alignment accuracy [6]. Another pairwise folding and alignment algorithm called *Consan* [5] first finds the confidently aligned base positions, referred as “pins”, and constrains the RNA alignment by fixing these positions. This is illustrated in Fig. 2(b). The set of pins  $\mathcal{P}$  is estimated using a pair-HMM, by looking for base positions  $(m, n)$  whose alignment probability  $P(x_m \Leftrightarrow y_n | \mathbf{X}, \mathbf{Y})$  exceeds some threshold  $p_h$ , which

<sup>2</sup>Two bases  $x_m$  and  $y_n$  are said to be co-incident if (i) they are either aligned to each other, or (ii) if  $x_m$  is inserted in the region that immediately follows  $x_{m_1}$  ( $m_1 < m$ ) which is aligned to  $y_n$ , or vice versa. [6].

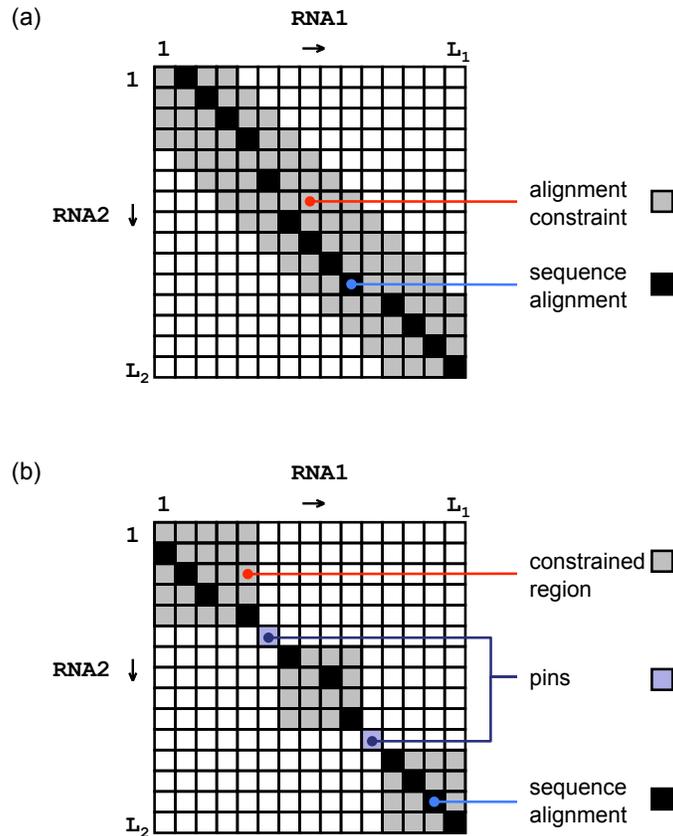


Fig. 2. (a) Alignment constraint in Dynalign [11]. The maximum distance between the aligned bases is restricted. (b) Consan [5] constrains the alignment space by fixing the “pins”, or confidently aligned base positions.

is close to unity. This set  $\mathcal{P}$  can be written as follows

$$\mathcal{P} = \left\{ (m, n) \mid P(x_m \Leftrightarrow y_n | \mathbf{X}, \mathbf{Y}) \geq p_h \right\}. \quad (3)$$

For every predicted pin  $(m, n) \in \mathcal{P}$ , the bases  $x_m$  and  $y_n$  are forced to be aligned to each other in the final alignment. While Dynalign [6] finds the set of *all possible* alignment positions, Consan [5] tries to find only a small set of alignment positions that *must* be included in the final alignment.

Although the previous alignment constraints [5], [6], [11] were mainly used to speed up Sankoff-style joint alignment and folding algorithms, similar ideas can be used to expedite dynamic programming alignment algorithms such as the CYK algorithm [3] for CMs and the *SCA (sequential component adjoining) algorithm* [16], [17] for profile-csHMMs. In the following section, we propose a new method for finding effective sequence alignment constraints that are especially useful for making these algorithms faster.

### III. ALIGNMENT CONSTRAINTS FOR RNA FAMILY-SPECIFIC MODELS

Let us assume that we have a reference RNA whose structure is known. This can be either the consensus sequence of an RNA family or simply a single RNA sequence. Also assume that we are given a target RNA with an unknown structure, which might be a putative member of the same family. We want to find the optimal alignment between these RNAs by considering both their sequence and structural similarities.

This *structural alignment* can be used for predicting the secondary structure of a new homologue [17], [20] or performing an RNA similarity search to identify new members in the same RNA family [3]. In order to find the structural alignment, we first construct a stochastic model (such as a profile-csHMM or a CM) that can closely represent the reference RNA. Then we use a dynamic programming alignment algorithm to find the best alignment between the reference RNA (represented by the constructed model) and the target RNA. Although the computational complexity of these algorithms is generally lower than that of Sankoff-style algorithms, it still ranges between  $O(L^3)$  and  $O(L^6)$  for a target RNA of length  $L$ .<sup>3</sup> This renders the dynamic programming algorithms impractical for aligning long RNAs or scanning a large database, and using effective alignment constraints can be greatly helpful in relieving this problem.

#### A. Motivation for estimating constraints based on predicted alignment positions

When we are interested in a specific RNA family, it will be more appropriate to establish the alignment constraints based on the member sequences in the given family. Therefore, it will be more desirable to use a *family-specific* model for finding the constraints, rather than using a *general* model that applies to all RNAs as in [5] and [6]. However, in many practical situations, we may not have enough number of sequences in the given family for reliably estimating the model parameters. As we can see in (2) and (3), the alignment constraints in Dynalign [6] and Consan [5] strongly depend on the estimated alignment probabilities. Although the alignment constraints used in these methods are expected to work well when we have a large number of training sequences, they are not suitable when only a handful of RNAs are available for training the model.

So, how can we find efficient alignment constraints for a family-specific model when we have only a limited number of sequences in the reference RNA family? In order to answer this question, let us consider the pair-HMMs shown in Fig. 3. Both pair-HMMs have three hidden states, ALN,  $l_X$ , and  $l_Y$ , for base alignment, base insertion in  $\mathbf{X}$  (RNA-1), and base insertion in  $\mathbf{Y}$  (RNA-2), respectively. The state ALN emits a pair of aligned bases  $x_m \in \mathbf{X}$  and  $y_n \in \mathbf{Y}$ . The insert state  $l_X$  emits an unaligned base  $x_m$  in  $\mathbf{X}$ , and similarly, the state  $l_Y$  emits an unaligned base  $y_n$  in  $\mathbf{Y}$ . These pair-HMMs can be used for finding a *sequence-based alignment*<sup>4</sup> between two RNAs, and for estimating the base alignment probabilities. Similar models have been used to find the alignment constraints (2) and (3) in Dynalign [6] and Consan [5], respectively. In this example, the transition probabilities of the pair-HMMs are shown along the arrows. We assume that the probability of entering a state in the beginning is identical for all three states. The emission probability  $e(x, y|\text{ALN})$  of a pair of aligned bases  $(x, y)$  is shown inside the box below the respective HMMs. Finally, the emission probability at an insert state is specified as follows

$$e(x|s) = \frac{1}{4}, \quad \forall x \in \{A, C, G, U\}, \quad \forall s \in \{l_X, l_Y\}. \quad (4)$$

Now, let us assume that we want to find the sequence-based alignment of the following RNAs:

$$\mathbf{X} = \text{AACUG}$$

$$\mathbf{Y} = \text{CUGAA},$$

<sup>3</sup>For RNAs without pseudoknots, the computational complexity of the alignment algorithm will be  $O(L^3)$ . The complexity for aligning pseudoknotted RNAs is at least  $O(L^4)$ , and it can become higher as the structure gets more complex.

<sup>4</sup>It is called a *sequence-based alignment*, since the alignment is obtained based on sequence similarity alone.

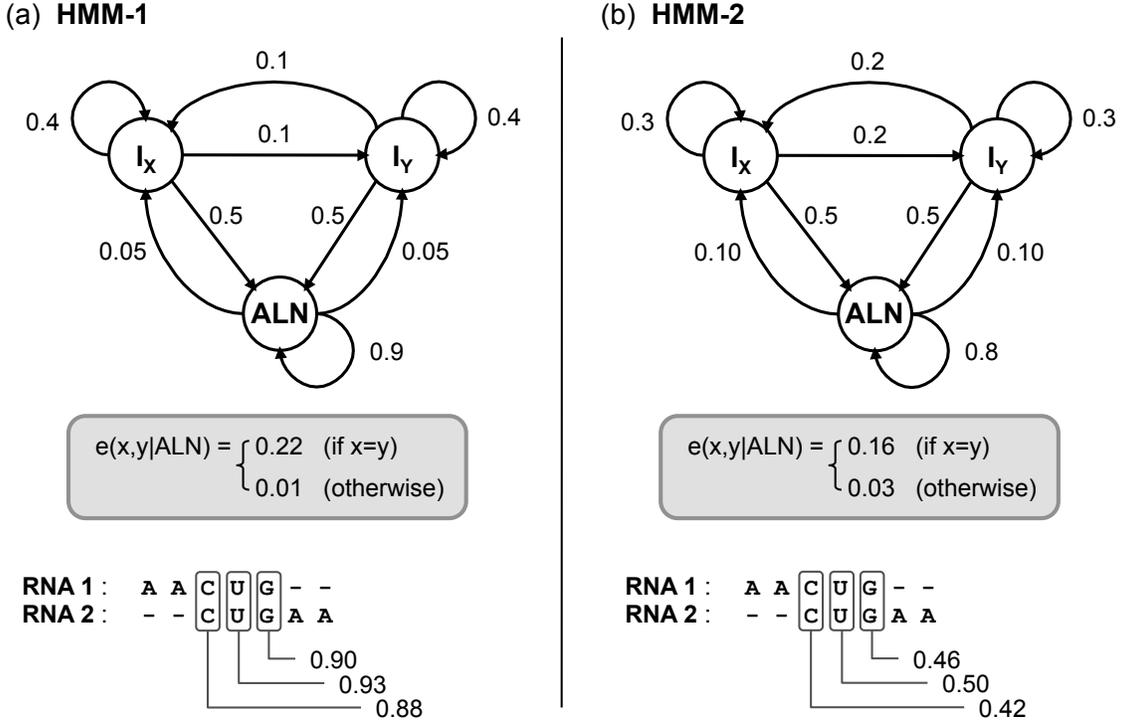


Fig. 3. Two pair-HMMs with slightly different parameters. Both pair-HMMs have three states, ALN,  $I_x$ , and  $I_y$ , which represent base alignment, base insertion in RNA-1, and base insertion in RNA-2, respectively.

using the pair-HMM shown in Fig. 3(a). Using the Viterbi algorithm, we can get the following alignment

$$\begin{aligned} \mathbf{X} : & A A C U G - - \\ \mathbf{Y} : & - - C U G A A. \end{aligned} \quad (5)$$

For each aligned pair  $(x_m, y_n)$ , we can compute the alignment probability  $P(x_m \Leftrightarrow y_n | \mathbf{X}, \mathbf{Y})$  using the *forward-backward algorithm* [21]. The estimated base alignment probabilities are shown in Fig. 3(a), below the RNA alignment. The estimated alignment probabilities are close to unity, indicating that we can be more or less confident about the predicted base alignments. Now, let us repeat this process using the pair-HMM shown in Fig. 3(b), which has slightly different parameters. As we can see in Fig. 3(b), HMM-2 finds the same alignment as HMM-1, but the estimated alignment probabilities are significantly different from the previous estimates. This example clearly shows that the estimation of the base alignment probability  $P(x_m \Leftrightarrow y_n | \mathbf{X}, \mathbf{Y})$  can be very sensitive to small changes in the model parameters. This implies that the alignment constraint in (3), which depends on  $P(x_m \Leftrightarrow y_n | \mathbf{X}, \mathbf{Y})$ , may not be reliable when we do not have enough training data to accurately estimate the HMM parameters. Compared to this, the alignment constraint in (2) might be more reliable, as the co-occurrence probability  $P(x_m \leftrightarrow y_n | \mathbf{X}, \mathbf{Y})$  used to estimate the constraint also includes the insertion probabilities. However, the predicted constraint will nevertheless depend on the parameters of the HMM to a considerable extent.

However, one thing we can notice by comparing Fig. 3(a) and Fig. 3(b) is that, despite the large difference in the estimated alignment probabilities, the resulting sequence alignments are identical. In fact, the alignment positions in an optimal sequence alignment are not very sensitive to small parameter

changes, and as a result, HMMs with reasonably similar parameters often yield almost identical alignment results. This motivates us to exploit the *aligned base positions*  $(x_m, y_n)$  for establishing the alignment constraints, instead of using the *base alignment probabilities*  $P(x_m \Leftrightarrow y_n | \mathbf{X}, \mathbf{Y})$ .

### B. Finding seed regions using a profile-HMM

Based on the previous observation, we propose a new method that utilizes the predicted alignment positions to find the alignment constraints. As before, let us denote the structured reference RNA as  $\mathbf{X}$  (RNA-1) and the unstructured target RNA as  $\mathbf{Y}$  (RNA-2). Ultimately, we want to find the *structural alignment* of these RNAs. However, since the dynamic programming algorithm for finding the structural alignment is computationally expensive, we want to come up with effective alignment constraints that can speed up the alignment.

For this purpose, we first build a profile-HMM [3] based on the reference RNA family. This model is used to find the *sequence-based alignment* between the reference RNA (represented by the profile-HMM) and the target RNA. Secondly, we identify the regions that consist of multiple consecutive base alignments, or base matches. Although a single base match may not be meaningful by itself, having a region of consecutive matches often indicates that the alignment in the given region is reasonably accurate. This is especially true for those matches that are located in the middle of the region. For example, we can see in both Fig. 3(a) and Fig. 3(b) that the alignment probability  $P(x_m \Leftrightarrow y_n | \mathbf{X}, \mathbf{Y})$  is largest for the alignment between  $x_4 = \text{U}$  and  $y_2 = \text{U}$ , which is located between the matches  $(x_3, y_1)$  and  $(x_5, y_3)$ . Therefore, we exclude the matches near the end and keep the remainder to obtain a set of reliable base alignments. The set of reliable contiguous matches is referred as the *seed region*. The procedure for finding the seed regions can be summarized as follows:

- 1) Find a sequence-based alignment between the RNAs.
- 2) Identify all regions, and longest such regions, that consist of consecutive matches. Let  $N_{match}$  be the number of consecutive matches in a given region. Keep only those regions with  $N_{match} \geq \Gamma$ .
- 3) In each region, exclude the first  $\gamma$  matches in the left end and the last  $\gamma$  matches in the right end.
- 4) The region that consists of the  $N_{match} - 2\gamma$  remaining matches is defined as a seed region.

The integer parameters  $\Gamma$  and  $\gamma$  ( $< \Gamma/2$ ) define the seed regions during this process. In general, using a larger  $\Gamma$  will identify a smaller number of seed regions, and a larger  $\gamma$  makes the seed regions contain fewer but more reliable base matches. Fig. 4(a) illustrates an example alignment with three seed regions.

### C. Constraints in a seed region

Assume that we have identified  $K$  seed regions according to the procedure described in Sec. III-B. Since the base alignments in these regions are relatively reliable, we keep the alignment space in these regions small. Let us consider the  $k$ -th seed region illustrated in Fig. 4(b). We denote the begin index and the end index of the  $k$ -th seed in RNA-1 as  $\sigma_1^\ell(k)$  and  $\sigma_1^r(k)$ , respectively.<sup>5</sup> Similarly, the begin and the

<sup>5</sup>The superscripts  $\ell$  and  $r$  in  $\sigma_1^\ell(k)$  and  $\sigma_1^r(k)$  stand for ‘left’ and ‘right’, respectively.

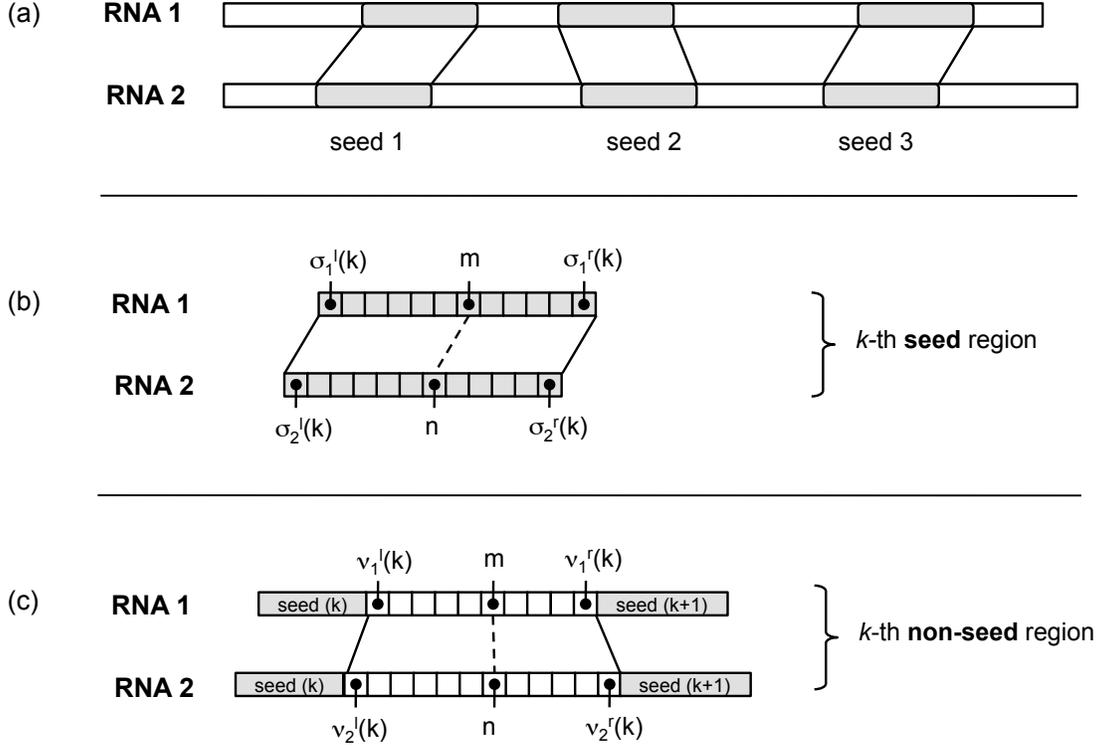


Fig. 4. Illustration of the proposed method. (a) Seed regions are identified from the sequence-based alignment. (b) Example of a seed region, which consists of consecutive base matches. (c) Example of a non-seed region.

end indices of the  $k$ -th seed in RNA-2 are denoted as  $\sigma_2^l(k)$  and  $\sigma_2^r(k)$ , respectively. Since a seed region consists of consecutive base matches, we have

$$\sigma_1^r(k) - \sigma_1^l(k) + 1 = \sigma_2^r(k) - \sigma_2^l(k) + 1 = L_k, \quad (6)$$

where  $L_k$  is the length of the  $k$ -th seed. For convenience, we define  $D_k$  as the position difference between the aligned bases in the  $k$ -th seed

$$\sigma_2^l(k) - \sigma_1^l(k) = \sigma_2^r(k) - \sigma_1^r(k) = D_k. \quad (7)$$

Based on the  $k$ -th seed, we define the set of allowed alignment positions  $(m, n)$  as follows

$$\mathcal{S}(k) = \left\{ (m, n) \mid \sigma_1^l(k) \leq m \leq \sigma_1^r(k), |n - m - D_k| \leq \Delta \right\}. \quad (8)$$

For a base  $x_m$  that is contained in the  $k$ -th seed region of RNA-1, the parameter  $\Delta$  restricts the distance between the base  $y_{n^*}$  ( $n^* = m + D_k$ ), to which  $x_m$  is aligned in the sequence-based alignment, and the base  $y_n$ , to which  $x_m$  will be aligned in the final structural alignment. As the base alignments in the seed regions are reliable,  $\Delta$  can be typically set to a small number. We find the set  $\mathcal{S}(k)$  for all  $k = 1, 2, \dots, K$ , and these sets will be combined later to establish the final alignment constraints.

#### D. Constraints in a non-seed region

The predicted base alignments in the non-seed regions are generally less reliable compared to those inside the seed regions. Therefore, we define different alignment constraints for the bases contained in

the non-seed regions, and make the constraints less stringent compared to (8). Let us consider the  $k$ -th non-seed region illustrated in Fig. 4(c). The begin and end indices of the  $k$ -th non-seed region in RNA-1 are denoted by  $\nu_1^\ell(k)$  and  $\nu_1^r(k)$ , respectively. Similarly,  $\nu_2^\ell(k)$  and  $\nu_2^r(k)$  respectively denote the begin and end indices of the corresponding non-seed region in RNA-2. Now, we define the set  $\mathcal{A}(k)$ , which contains (i) all aligned base positions  $(m, n)$  in the  $k$ -th non-seed region as well as (ii) the first and last positions  $(\nu_1^\ell(k), \nu_2^\ell(k))$  and  $(\nu_1^r(k), \nu_2^r(k))$  in this region.

$$\mathcal{A}(k) = \left\{ (m, n) \mid \nu_1^\ell(k) \leq m \leq \nu_1^r(k), \nu_2^\ell(k) \leq n \leq \nu_2^r(k), x_m \Leftrightarrow y_n \right\} \cup \left\{ (\nu_1^\ell(k), \nu_2^\ell(k)), (\nu_1^r(k), \nu_2^r(k)) \right\}. \quad (9)$$

In practice, it is possible that there may be no aligned bases  $x_m \Leftrightarrow y_n$  in the given non-seed region. Including the terminal positions  $(\nu_1^\ell(k), \nu_2^\ell(k))$  and  $(\nu_1^r(k), \nu_2^r(k))$  of the  $k$ -th non-seed region in  $\mathcal{A}(k)$  ensures that the set  $\mathcal{A}(k)$  will never be empty. For the position pairs  $(m, n) \in \mathcal{A}(k)$ , we estimate the range of the position difference  $(n - m)$  as follows

$$\Delta_{min}(k) = \min_{(m, n) \in \mathcal{A}(k)} (n - m) \quad (10)$$

$$\Delta_{max}(k) = \max_{(m, n) \in \mathcal{A}(k)} (n - m). \quad (11)$$

Based on these values, we define the following set

$$\mathcal{N}(k) = \left\{ (m, n) \mid \nu_1^\ell(k) \leq m \leq \nu_1^r(k), m + \Delta_{min}(k) - \Delta \leq n \leq m + \Delta_{max}(k) + \Delta \right\}, \quad (12)$$

which contains the alignable base positions  $(m, n)$  in the  $k$ -th non-seed region. Note that we use the same  $\Delta$  in (8) and (12). Therefore, a larger  $\Delta$  will relax the alignment constraints for both seed and non-seed regions, and a smaller  $\Delta$  will make both constraints more stringent.

#### E. Overall alignment constraints

In Sec. III-C and Sec. III-D, we defined the alignment constraints in the seed regions as well as the constraints in the non-seed regions. Finally, we combine (8) and (12) to obtain the overall alignment constraints  $\mathcal{C}$  as follows

$$\mathcal{C} = \left[ \bigcup_{\forall k_1} \mathcal{S}(k_1) \right] \cup \left[ \bigcup_{\forall k_2} \mathcal{N}(k_2) \right]. \quad (13)$$

This set  $\mathcal{C}$  can be used to constrain the alignment space of the dynamic programming algorithm for finding the structural alignment of the given RNAs. When finding the RNA alignment, we allow a base  $x_m$  in RNA-1 to be aligned to a base  $y_n$  in RNA-2 only if the pair  $(m, n)$  is included in this set  $\mathcal{C}$ .

## IV. EXPERIMENTAL RESULTS

To demonstrate the effectiveness of the proposed method, we applied the new alignment constraints to the profile-csHMM based structural alignment method [17]. In the following, we provide a brief explanation about the experimental set-up and present the experimental results.

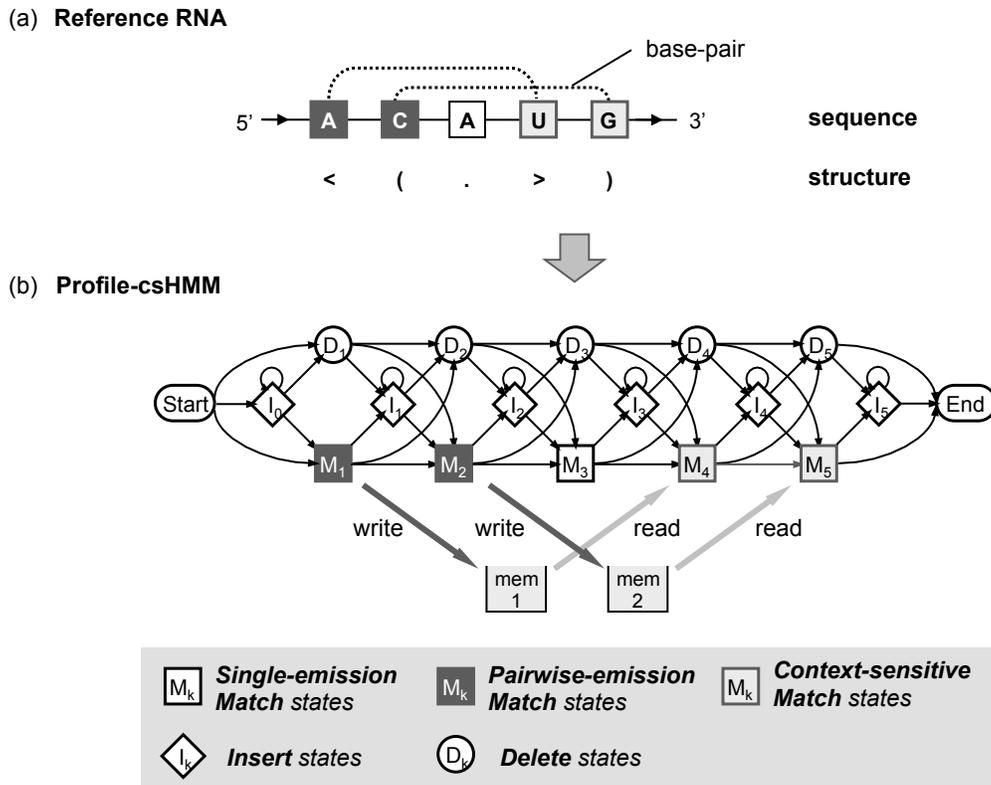


Fig. 5. Constructing a profile-csHMM. (a) A reference RNA sequence with a known secondary structure. (b) The profile-csHMM that represents the given reference RNA.

### A. Profile-csHMM based structural alignment

Profile-csHMMs are a subclass of context-sensitive HMMs [22] that are especially useful for representing RNA sequence profiles and their secondary structure. In principle, profile-csHMMs can represent RNA secondary structures with any kind of base-pairs [16], [17]. As a result, profile-csHMMs can also be used for aligning and predicting the structure of RNAs that contain pseudoknots, which cannot be done using the widely used SCFGs (or CMs). The profile-csHMM based structural alignment algorithm proposed in [17] proceeds as follows. In the first place, a profile-csHMM is constructed based on a reference RNA sequence with a known structure. In [17], a single reference RNA was used to build the model. This can be used for performing a single RNA homology search, similar to the CM-based search proposed in [23]. Figure 5 illustrates an example, where a profile-csHMM is constructed based on a reference RNA that has two crossing base-pairs. Obviously, we do not have enough training sequences to accurately estimate the model parameters in this case, hence the parameters of the profile-csHMM are chosen according to the scoring scheme proposed by Gorodkin and co-workers [24]. These scores can be viewed as normalized log-probabilities for observing base substitutions or gaps (insertions and deletions) in homologous RNAs. They have been used in a number of RNA alignment algorithms [20], [24], yielding accurate alignment results. The constructed profile-csHMM can then be used for finding the optimal structural alignment between the reference RNA and an unstructured target RNA, computing their alignment score, and

predicting the secondary structure of the target RNA. A dynamic programming alignment algorithm called the *SCA (sequential component adjoining) algorithm* can be used for this purpose.

### B. Estimating the alignment constraints

In order to estimate the alignment constraints for expediting the profile-csHMM alignment algorithm (or the SCA algorithm), we construct a profile-HMM based on the same reference RNA. Note that unlike the profile-csHMM, the traditional profile-HMM reflects only the sequence characteristics of the reference RNA. Similar to the parameterization of the profile-csHMM described in Sec. IV-A, the parameters of the profile-HMM are also specified according to the scores in [24]. The resulting profile-HMM is used to estimate the sequence alignment constraints as we elaborated in Sec. III. We use the estimated constraints to restrict the alignment space of the structural RNA alignment to reduce the overall computational load.

### C. Choosing the parameters for constraint estimation

Now, one practical question is how we should choose the parameters  $\Gamma$ ,  $\gamma$ , and  $\Delta$  that are used to estimate the alignment constraints in Sec. III-C and Sec. III-D. Ideally, the predicted alignment constraints should minimize the alignment space without affecting the quality of the final structural alignment. Since the alignment constraints proposed in Sec. III are derived from the predicted seed regions, the alignment accuracy in these regions has a crucial impact on the accuracy of the proposed approach. For this reason, we estimated the average base alignment accuracy in the seed regions for the 5S rRNA and tRNA families in the Rfam database (version 8.1) [25]. We used the RNAs in the *seed alignment* of the respective family, as they have a relatively reliable secondary structure annotation. For each RNA family, we first chose a reference RNA among its members, and constructed a profile-HMM based on the chosen RNA. Then we aligned the remaining members to the reference RNA using the profile-HMM. For every sequence alignment, the predicted alignment positions have been compared to the correct positions in the database to estimate the alignment error rate. In order to get a reliable estimate, we repeated these experiments by using every member as the reference RNA. This resulted in 1,182,656 pairwise alignments for tRNAs and 345,156 alignments for 5S rRNAs.

Let us assume that the profile-HMM predicted that  $x_m$  in the reference RNA should be aligned to  $y_n$  in the target RNA. We want to estimate the probability of error for this prediction as a function of the following parameters:

- 1)  $L_{aln}$ : the number of consecutive base matches in the region containing the alignment  $(m, n)$ ,
- 2)  $d_{min}$ : the minimum distance between the given base alignment  $(m, n)$  and the terminal alignment positions  $(m_1, n_1)$  and  $(m_2, n_2)$ . See Fig. 6 for illustration.

Consider the example illustrated in Fig. 6. The number of consecutive base matches, or the *length*, of the region containing  $(m, n)$  is

$$L_{aln} = m_2 - m_1 + 1 = n_2 - n_1 + 1, \quad (14)$$

and the minimum distance  $d_{min}$  is defined as

$$d_{min} = \min \left( (m - m_1), (m_2 - m) \right) = \min \left( (n - n_1), (n_2 - n) \right). \quad (15)$$

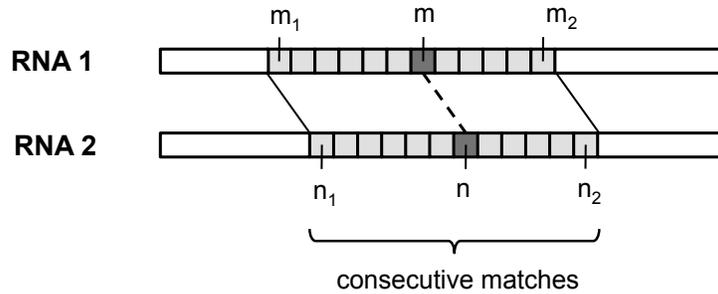


Fig. 6. A region that consists of consecutive base matches.

Based on these definitions, we define

$$P_e(L, d) = P\left(x_m \not\leftrightarrow y_n \text{ in the trusted alignment} \mid L_{aln} \geq L, d_{min} \geq d\right). \quad (16)$$

This is the probability that the predicted base alignment  $(m, n)$  between  $x_m$  and  $y_n$  will be incorrect, given that

- 1) the length  $L_{aln}$  of the alignment region containing  $(m, n)$  is at least  $L$ , and
- 2) there are at least  $d$  matches in the left-hand side of  $(m, n)$  as well as in the right-hand side.

Fig. 7(a) shows the contour plot of the misalignment probability  $P_e(L, d)$  for 5S rRNAs, where the  $x$ -axis is for  $L$  and the  $y$ -axis is for  $d$ . On top of each contour curve, we show the corresponding misalignment probability  $P_e(L, d)$  for the points  $(L, d)$  on the given curve. Darker shaded regions correspond to higher  $P_e(L, d)$  and lighter shaded regions correspond to lower  $P_e(L, d)$ . The diagonal line representing  $L = 2d + 1$  is shown in the plot as a reference. Note that, by definition, we have  $L \geq 2d + 1$ . Therefore, for any  $(L, d)$  such that  $L < 2d + 1$ , which corresponds to the region above the diagonal line, we will have  $P_e(L, d) = P_e(L, \frac{L-1}{2})$ .

As we would expect, the misalignment probability  $P_e(L, d)$  becomes smaller as  $L$  and  $d$  get larger. Fig. 7(b) shows the misalignment probability  $P_e(L, d)$  of 5S rRNAs for  $L = 2d + 1$ . The pairwise alignments have been divided into different groups based on the *percentage identity* (or *percent sequence similarity*) of the aligned RNAs, and the alignment error probability  $P_e(L, d)$  has been computed for the respective groups. As we can see in Fig. 7(b), the error probability is generally lower for RNAs with higher percentage identity. This is expected, since the seed regions are predicted from a sequence-based alignment, which will be more accurate if the RNAs have higher sequence similarity. Fig. 7(c) and Fig. 7(d) show the misalignment probability  $P_e(L, d)$  for tRNAs, which have similar trends.

We also computed  $P_e(L, d)$  for RNAs with 60%~100% identity for different values of  $(L, d)$ . This is summarized in Table I. For example, the alignment error probability  $P_e(L, d)$  for 5S rRNAs is 1.81% for  $(L, d) = (9, 4)$  and 0.79% for  $(L, d) = (15, 7)$ . This implies that if we choose  $\Gamma = 9$  and  $\gamma = 4$  when finding the seed regions (see Sec. III-B), more than 98% of the alignments in the predicted seed regions will be correct. In our experiments, we observed that the misaligned bases were typically located within 1~2 base positions from the correct ones. This implies that if we let  $\Delta = 1$  or  $\Delta = 2$ , most of the correct base alignments  $(m, n)$  will be included in the constrained alignment space  $\mathcal{S}(k)$  in (8). Therefore, imposing

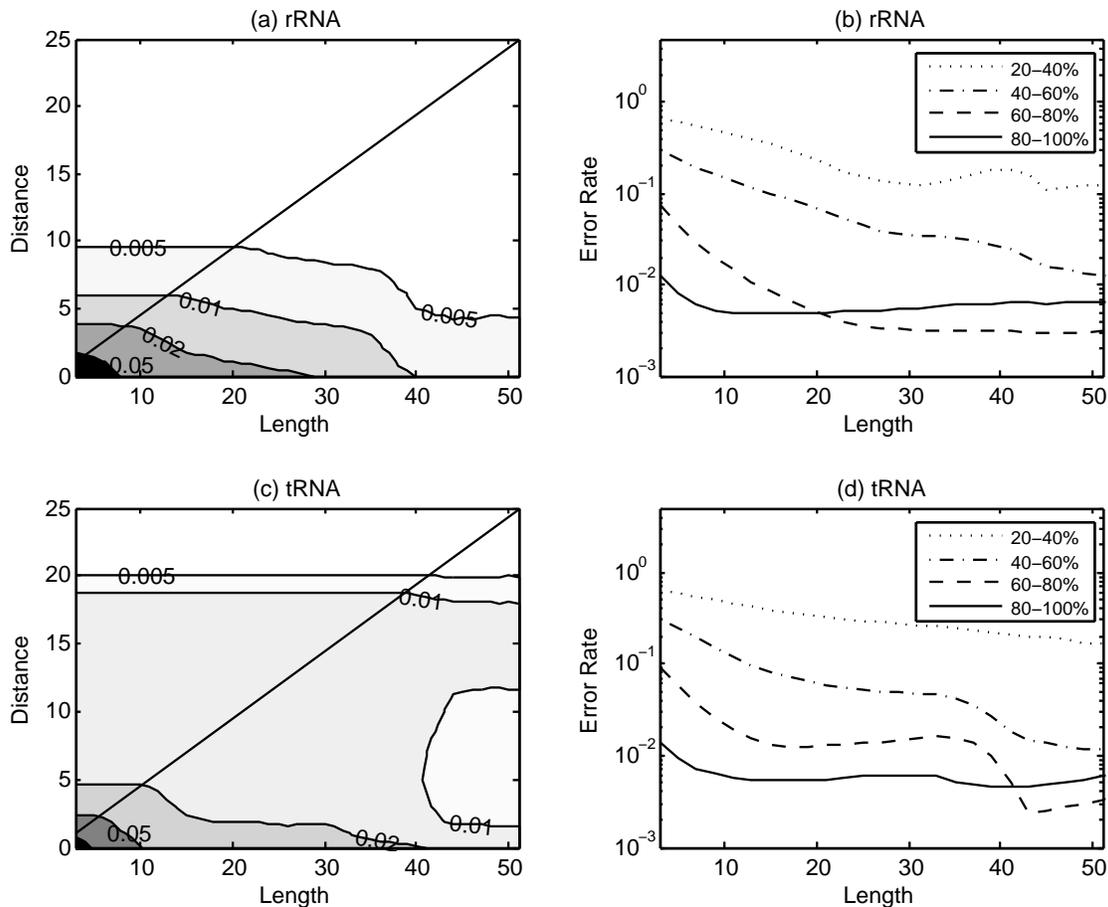


Fig. 7. Alignment error probability  $P_e(L, d)$  in the seed regions. (a) Contour plot of  $P_e(L, d)$  for 5S rRNAs. Darker shade corresponds to higher  $P_e(L, d)$  and lighter shade corresponds to lower  $P_e(L, d)$ . The misalignment probability  $P_e(L, d)$  on the level curves are shown on top of the contours. (b) Misalignment probability for 5S rRNAs with different percentage identities. (c) Contour plot of  $P_e(L, d)$  for tRNAs. (d) Misalignment probability for tRNAs with different percentage identities.

TABLE I

MISALIGNMENT PROBABILITY IN THE SEED REGIONS OF RNAs WITH 60%~100% IDENTITY. THE PROBABILITY  $P_e(L, d)$  HAS BEEN ESTIMATED FOR FIVE DIFFERENT VALUES OF  $(L, d)$ .

	$(L, d)$				
	(7,3)	(9,4)	(11,5)	(13,6)	(15,7)
	ERROR (%)				
5S RIBOSOMAL RNA	2.64	1.81	1.32	1.00	0.79
TRANSFER RNA	3.60	2.46	1.80	1.44	1.28

these constraints will not degrade the accuracy of the final structural alignment.

#### D. RNA structural alignment with the proposed constraints

As mentioned earlier, we applied the proposed constraints to the profile-csHMM based structural alignment algorithm [17]. For our experiments, we chose six RNA families from the Rfam<sup>6</sup> database [25]. Among the six families, two families, tRNAs and 5S rRNAs, do not contain pseudoknots in their

<sup>6</sup>The Flavi.pk3 RNAs were obtained from Rfam 7.0, which are now part of the Flavi.CRE family in Rfam 8.1. For other families, we obtained the RNAs from Rfam 8.1.

TABLE II  
BASIC PROPERTIES OF THE RNA FAMILIES USED IN THE EXPERIMENTS.

	# OF SEED SEQUENCES	AVERAGE LENGTH	AVERAGE PERCENTAGE IDENTITY
TRANSFER RNA	1088	72.7	45
5S RIBOSOMAL RNA	602	116.8	61
CORONA_PK3	14	62.5	70
HDV_RIBOZYME	15	88.8	95
TOMBUS_3_IV	18	64.5	94
FLAVI_PK3	14	95.4	69

secondary structures, while the other four families, Corona\_pk3, HDV\_ribozyme, Tombus\_3.IV, and Flavi\_pk3, contain pseudoknots. The basic properties of these RNA families, such as the number of RNAs in the Rfam seed alignment, the average length of the member sequences, and their average percentage (sequence) identity, are shown in Table II. For each RNA family, we performed the following experiment:

- 1) Choose a reference RNA from the seed alignment.
- 2) Construct a profile-HMM and a profile-csHMM based on the reference RNA.
- 3) Choose a different target RNA from the seed alignment.
- 4) Estimate the alignment constraint using the profile-HMM.
- 5) Apply the alignment constraint and find the structural alignment using the profile-csHMM.
- 6) Repeat STEP-3 to STEP-5 for different target RNAs.
- 7) Repeat STEP-1 to STEP-6 for different reference RNAs.

In order to measure the quality of the structural alignment, we predicted the secondary structure of the target RNA based on the structural alignment, and compared it to the trusted structure in the Rfam database. Then we counted the number of correctly predicted base-pairs (TP; true-positives), the number of incorrectly predicted base-pairs (FP; false-positives), and the number of true base-pairs that could not be predicted (FN; false-negatives). Based on these numbers, we estimated the *sensitivity* (SN) and the *positive predictive value* (PPV) as follows

$$SN = \frac{TP}{TP + FN}, \quad PPV = \frac{TP}{TP + FP}. \quad (17)$$

The sensitivity is defined as the fraction of base-pairs in the trusted structure that could be predicted by the algorithm, and the positive predictive value is defined as the fraction of predicted base-pairs that were correct.

We first tested the performance of the proposed approach on RNA families that do not contain pseudoknots. In order to compare the effectiveness of different alignment constraints, we repeated the above experiment for the following methods:

- 1) Profile-csHMM + proposed alignment constraint (referred as "PROPOSED")
- 2) Profile-csHMM + M-constraint (referred as "M-CONSTRAINT")
- 3) Profile-csHMM (original implementation in [17]; referred as "ORIGINAL")

TABLE III  
AVERAGE SENSITIVITY (SN), POSITIVE PREDICTIVE VALUE (PPV), AND ALIGNMENT TIME FOR RNA FAMILIES THAT DO NOT CONTAIN PSEUDOKNOTS.

	PROFILE-CSHMM								
	M-CONSTRAINT			ORIGINAL			PROPOSED		
	SN (%)	PPV (%)	TIME (sec)	SN (%)	PPV (%)	TIME (sec)	SN (%)	PPV (%)	TIME (sec)
TRANSFER RNA	94.2	95.8	0.0739	94.1	96.0	0.0139	93.6	96.2	0.0108
5S RIBOSOMAL RNA	94.8	96.3	0.0676	95.1	97.0	0.0024	95.9	98.5	0.0010

Table III summarizes the average sensitivity (SN), positive predictive value (PPV), and alignment time<sup>7</sup> for using different alignment constraints with the profile-csHMM based structural alignment method. These results have been obtained from one thousand structural alignments of distinct pairs of RNAs that were chosen from the seed alignment of the respective RNA families. For these experiments, we used  $(\Gamma, \gamma, \Delta)=(9, 4, 2)$  for tRNAs and  $(\Gamma, \gamma, \Delta)=(9, 4, 0)$  for 5S rRNAs. These parameters were chosen based on the analysis in Sec. IV-C.<sup>8</sup> For the M-constraint defined in (1), we used  $M = 7$  as in [6]. As we can see in Table III, all three methods were able to achieve accurate alignment results that were comparable to each other. However, adopting the proposed alignment constraint improved the average alignment speed significantly, which was around 7~68 times faster compared to the fixed M-constraint, and up to 2.4 times faster compared to the original implementation in [17] that uses a simple heuristic.

In order to test the performance of the proposed method on RNA pseudoknots, we carried out similar experiments using four pseudoknotted RNA families, Corona\_pk3, HDV\_ribozyme, Tombus\_3\_IV and Flavi\_pk3. For these experiments, we used  $(\Gamma, \gamma, \Delta)=(9, 4, 0)$  and  $M = 3$  for all four RNA families. In addition to evaluating the performance of the profile-csHMM method for these families, we evaluated the performance of the PSTAG-based method [20] for comparison. The PSTAG-based structural alignment method is a state-of-the-art pairwise RNA alignment method that uses *pair stochastic tree adjoining grammars (PSTAGs)*. PSTAGs can be used for aligning many known pseudoknots, though not all of them. To the best of our knowledge, the PSTAG-based alignment method [20] is the only grammar-based method that can be used for finding the structural alignment of pseudoknotted RNAs, except for the profile-csHMM method. Table IV shows the average sensitivity and positive predictive value of the different alignment methods.<sup>9</sup> As we can see in this table, all four methods could achieve high sensitivity and PPV for the Corona\_pk3, HDV\_ribozyme, and Tombus\_3\_IV RNA families. The Flavi\_pk3 RNAs could not be aligned using PSTAGs, as they have a more complex secondary structure compared to other RNA families. Unlike PSTAGs, profile-csHMMs can handle RNAs with any kind of base-pairs, hence they could be used for aligning Flavi\_pk3 RNAs as well.<sup>10</sup> Table IV shows that all three profile-csHMM based

<sup>7</sup>The CPU time for finding the alignment has been measured on a MacPro with two 2.8GHz quad-core Intel Xeon processors and 4GB memory.

<sup>8</sup>In general, there will be a trade-off between alignment accuracy and runtime. These parameters have been used as they provide a good balance between these two measures. Further discussion on this trade-off can be found at the end of Sec. IV-D.

<sup>9</sup>The sensitivity and the positive predictive value of the PSTAG-based method have been obtained from [20] based on the same test set.

<sup>10</sup>The current implementation can handle any RNAs in the Rivas&Eddy class [26], which includes nearly all known pseudoknots. We can also handle the RNAs outside the R&E class by incorporating additional *adjoining rules*. See [19] for further discussions on adjoining rules and the descriptive capability of profile-csHMMs.

TABLE IV  
AVERAGE SENSITIVITY (SN) AND POSITIVE PREDICTIVE VALUE (PPV) FOR RNA FAMILIES WITH PSEUDOKNOTS.

	PROFILE-CSHMM						PSTAG	
	M-CONSTRAINT		ORIGINAL		PROPOSED		SN (%)	PPV (%)
	SN (%)	PPV (%)	SN (%)	PPV (%)	SN (%)	PPV (%)		
CORONA_PK3	95.5	95.7	95.7	96.5	94.8	96.0	94.6	95.5
HDV_RIBOZYME	94.5	95.1	94.5	95.3	94.2	95.9	94.1	95.6
TOMBUS_3_IV	95.9	96.4	95.9	96.4	96.8	97.4	97.4	97.4
FLAVI_PK3	94.6	96.5	94.5	96.4	94.5	96.8	N/A	N/A

TABLE V  
AVERAGE CPU TIME FOR FINDING THE STRUCTURAL ALIGNMENT OF RNAs CONTAINING PSEUDOKNOTS.

	PROFILE-CSHMM			PSTAG
	M-CONSTRAINT	ORIGINAL	PROPOSED	TIME (SEC)
	TIME (SEC)	TIME (SEC)	TIME (SEC)	
CORONA_PK3	9.37	0.71	0.23	19.65
HDV_RIBOZYME	10.30	1.03	0.13	158.77
TOMBUS_3_IV	6.99	0.35	0.07	193.06
FLAVI_PK3	13.31	3.96	0.35	N/A

approaches yielded accurate alignment results for Flavi\_pk3 RNAs. By comparing the performance of the profile-csHMM method with different constraints, we can note that incorporating the proposed alignment constraint virtually did not affect the alignment accuracy.

As we can see in Table V, the proposed sequence alignment constraint was able to significantly improve the alignment speed also for pseudoknotted RNAs. In fact, by comparing the results in Table III and Table V, we can observe that the overall computational gain becomes even larger for RNAs with more complicated secondary structures. The new constraint made the alignment speed around 40~100 times faster compared to the fixed M-constraint (using  $M = 3$ ), and around 3~11 times faster compared to the original implementation [17], at a comparable prediction accuracy. We can also note that the PSTAG-based alignment takes considerably longer than the profile-csHMM based alignment. The large difference in alignment speed is mainly due to the fact that the PSTAG algorithm [20] does not incorporate any constraint to restrict the alignment space.

It would be also interesting to see how the parameters used for predicting the alignment constraint would affect the overall performance. For this purpose, we repeated the previous experiment using different values of  $\Gamma$  and  $\gamma$ . Three pairs of  $(\Gamma, \gamma)$  were chosen based on the experimental results shown in Fig. 7, such that the average misalignment probability does not exceed 5% for both tRNAs and 5S rRNAs. We used  $\Gamma = 2\gamma + 1$  in all three cases, such that the minimum length of the seed region is one. Note that if there are regions with more than  $\Gamma$  consecutive matches in the sequence-based alignment, the lengths of the corresponding seed regions will be longer than this minimum. In all three experiments, the parameter  $\Delta$  was set to zero. Table VI shows the sensitivity, PPV, and alignment time for different pairs of  $(\Gamma, \gamma)$ . In general, small  $\Gamma$  and  $\gamma$  tend to increase the fraction of bases included in the seed regions, thereby reducing the overall alignment space. As a consequence, the alignment time becomes smaller as we can see in Table VI. However, if these values are made too small, the resulting alignment space can become too restricted, hence degrading the alignment accuracy. This phenomenon could be

TABLE VI  
PERFORMANCE OF THE PROPOSED APPROACH FOR DIFFERENT PARAMETERS.

	PROFILE-CSHMM (PROPOSED)								
	$\Gamma = 7, \gamma = 3$			$\Gamma = 9, \gamma = 4$			$\Gamma = 11, \gamma = 5$		
	SN (%)	PPV (%)	TIME (SEC)	SN (%)	PPV (%)	TIME (SEC)	SN (%)	PPV (%)	TIME (SEC)
CORONA_PK3	92.9	94.4	0.139	94.8	96.0	0.232	95.3	96.2	0.278
HDV_RIBOZYME	93.2	95.7	0.131	94.2	95.9	0.133	94.5	95.5	0.147
TOMBUS_3_IV	96.5	97.1	0.065	96.8	97.4	0.068	96.6	97.3	0.069
FLAVI_PK3	94.8	97.2	0.329	94.5	96.8	0.351	94.5	96.8	0.362

observed when aligning the Corona\_pk3 RNAs with  $\Gamma = 7$  and  $\gamma = 3$ .

In Table III and Table V, we have shown that the proposed alignment constraint can significantly reduce the average computational requirement for finding the RNA structural alignments. Since the proposed method estimates the constraint based on the sequence alignment of the given RNAs, the actual reduction in complexity will depend on the degree of sequence similarity between the RNAs. Suppose we have a reference RNA of length  $L_r$  and a target RNA of length  $L_t$ . In the best case, when these RNAs are perfectly aligned, the overall computational cost will be dominated by the constraint estimation step, hence the resulting complexity will be  $O(L_r L_t)$ . In the worst case, the complexity will be identical to that of an unconstrained profile-csHMM alignment, which is  $O(L_r L_t^3)$  for RNAs without pseudoknots,  $O(L_r L_t^4)$  for typical RNA pseudoknots (including Corona\_pk3, HDV\_ribozyme, Tombus\_3\_IV, and Flavi\_pk3 used in our experiments), and  $O(L_r L_t^6)$  for RNAs with the most complicated secondary structure in the Rivas&Eddy class [26]. In general, the maximum distance between the alignable bases will be limited by the constraint (12). If we define  $D$  as

$$D = 2\Delta + \max_k \left( \Delta_{max}(k) - \Delta_{min}(k) \right), \quad (18)$$

the computational complexity of the profile-csHMM alignment method with the proposed constraint will be  $O(L_r L_t + L_r D^3)$  for RNAs that do not contain pseudoknots. For pseudoknotted RNAs in the Rivas&Eddy class, the complexity will range between  $O(L_r L_t + L_r D^4)$  and  $O(L_r L_t + L_r D^6)$ .

## V. CONCLUDING REMARKS

In this paper, we proposed a new method for finding an effective alignment constraint for fast and accurate structural alignment of RNAs. The proposed method is especially useful for accelerating the dynamic programming alignment algorithm of family-specific models, such as the profile-csHMMs or CMs. The alignment constraint proposed in this paper is not very sensitive to small parameter changes in the model that is used to predict the constraint. Therefore, it can be especially useful when we do not have enough sequences in the reference RNA family for training the model. We applied the new constraint to the profile-csHMM based structural alignment method [17], and evaluated its performance using several RNA families containing pseudoknots. Experimental results showed that the proposed alignment constraint could significantly reduce the alignment time without any loss of alignment accuracy. Although we have mainly focused on incorporating the proposed constraint into the profile-csHMM based method, these constraints can certainly be used to expedite other alignment methods based on CMs [3], [23] or PSTAGs [20].

## REFERENCES

- [1] S. R. Eddy, "Non-coding RNA genes and the modern RNA world", *Nature Reviews Genetics*, vol. 2, pp. 919-929, 2001.
- [2] B.-J. Yoon and P. P. Vaidyanathan, "Computational identification and analysis of noncoding RNAs - Unearthing the buried treasures in the genome", *IEEE Signal Processing Magazine*, vol. 24, no. 1, pp. 64-74, Jan. 2007.
- [3] R. Durbin, S. Eddy, A. Krogh, and G. Mitchison, *Biological Sequence Analysis*, Cambridge, UK: Cambridge University Press, 1998.
- [4] D. Sankoff, "Simultaneous solution of the RNA folding, alignment, and protosequence problems," *SIAM Journal on Applied Mathematics*, vol. 45, pp. 810-825, 1985.
- [5] R. D. Dowell and S. R. Eddy, "Efficient pairwise RNA structure prediction and alignment using sequence alignment constraints," *BMC Bioinformatics*, 7:400, 2006.
- [6] A. O. Harmanci, G. Sharma, and D. H. Mathews, "Efficient pairwise RNA structure prediction using probabilistic alignment constraints in Dynalign," *BMC Bioinformatics*, 8:130, 2007.
- [7] J. H. Havgaard, R. B. Lyngsø, G. D. Stormo, and J. Gorodkin, "Pairwise local structural alignment of RNA sequences with sequence similarity less than 40%," *Bioinformatics*, vol. 21, pp. 1815-1824, 2005.
- [8] J. H. Havgaard, E. Torarinsson, and J. Gorodkin, "Fast pairwise structural RNA alignments by pruning of the dynamical programming matrix," *PLoS Computational Biology*, 3(10):e193, 2007.
- [9] I. Holmes, "Accelerated probabilistic inference of RNA structure evolution," *BMC Bioinformatics*, 6:73, 2005.
- [10] D. H. Mathews, "Predicting a set of minimal free energy RNA secondary structures common to two sequences," *Bioinformatics*, vol. 21, no. 10, pp. 2246-2253, 2005.
- [11] A. V. Uzilov, J. M. Keegan, D. H. Mathews, "Detection of non-coding RNAs on the basis of predicted secondary structure formation free energy change," *BMC Bioinformatics*, 7: 173, 2006.
- [12] Z. Weinberg and W. L. Ruzzo, "Exploiting conserved structure for faster annotation of non-coding RNAs without loss of accuracy," *Bioinformatics*, vol. 20, Suppl.1:i334-41, 2004.
- [13] Z. Weinberg and W. L. Ruzzo, "Sequence-based heuristics for faster annotation of non-coding RNA families," *Bioinformatics*, vol. 22, no. 1, pp. 35-39, Jan. 2006.
- [14] B.-J. Yoon and P. P. Vaidyanathan, "Fast search of sequences with complex symbol correlations using profile context-sensitive HMMs and pre-screening filters", *Proc. 32nd International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Honolulu, Hawaii, Apr. 2007.
- [15] B.-J. Yoon and P. P. Vaidyanathan, "Fast structural similarity search of noncoding RNAs based on matched filtering of stem patterns", *Proc. 41st Asilomar Conference on Signals, Systems, and Computers*, Monterey, CA, Nov. 2007.
- [16] B.-J. Yoon and P. P. Vaidyanathan, "Profile context-sensitive HMMs for probabilistic modeling of sequences with complex correlations", *Proc. 31st International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Toulouse, May 2006.
- [17] B.-J. Yoon and P. P. Vaidyanathan, "Structural alignment of RNAs using profile-csHMMs and its application to RNA homology search: Overview and new results," *IEEE Transactions on Automatic Control (Joint Special Issue on Systems Biology with IEEE Transactions on Circuits and Systems: Part-I)*, vol. 53, pp. 10-25, Jan. 2008.
- [18] E. P. Nawrocki and S. R. Eddy, "Query-dependent banding (QDB) for faster RNA similarity searches," *PLoS Computational Biology*, 3(3): e56, 2007.
- [19] B.-J. Yoon and P. P. Vaidyanathan, "Fast Structural Alignment of RNAs by Optimizing the Adjoining Order of Profile-csHMMs," *IEEE Journal of Selected Topics in Signal Processing*, accepted.
- [20] H. Matsui, K. Sato, and Y. Sakakibara, "Pair stochastic tree adjoining grammars for aligning and predicting pseudoknot RNA structures", *Bioinformatics*, vol. 21, pp. 2611-2617, 2005.
- [21] L. R. Rabiner, "A tutorial on hidden Markov models and selected applications in speech recognition," *Proceedings of the IEEE*, vol. 77, pp. 257-286, 1989.
- [22] B.-J. Yoon and P. P. Vaidyanathan, "Context-sensitive hidden Markov models for modeling long-range dependencies in symbol sequences", *IEEE Transactions on Signal Processing*, vol. 54, pp. 4169-4184, Nov. 2006.
- [23] R. J. Klein and S. R. Eddy, "RSEARCH: Finding homologs of single structured RNA sequences," *BMC Bioinformatics*, vol. 4, 44, 2003.
- [24] J. Gorodkin, L. J. Heyer, and G. D. Stormo, "Finding the most significant common sequence and structure motifs in a set of RNA sequences", *Nucleic Acids Research*, vol. 25, pp. 3724-3732, 1997.
- [25] S. Griffiths-Jones, S. Moxon, M. Marshall, A. Khanna, S. R. Eddy and A. Bateman, "Rfam: annotating non-coding RNAs in complete genomes," *Nucleic Acids Res.*, vol. 33, pp. D121-D124, 2005.
- [26] E. Rivas and S. R. Eddy, "The language of RNA: A formal grammar that includes pseudoknots," *Bioinformatics*, vol. 16, pp. 334-340, 2000.