

RESEARCH

Open Access

# Hierarchical folding of multiple sequence alignments for the prediction of structures and RNA-RNA interactions

Stefan E Seemann<sup>†1</sup>, Andreas S Richter<sup>†2</sup>, Jan Gorodkin<sup>1</sup> and Rolf Backofen<sup>\*2</sup>

## Abstract

**Background:** Many regulatory non-coding RNAs (ncRNAs) function through complementary binding with mRNAs or other ncRNAs, e.g., microRNAs, snoRNAs and bacterial sRNAs. Predicting these RNA interactions is essential for functional studies of putative ncRNAs or for the design of artificial RNAs. Many ncRNAs show clear signs of undergoing compensating base changes over evolutionary time. Here, we postulate that a non-negligible part of the existing RNA-RNA interactions contain preserved but covarying patterns of interactions.

**Methods:** We present a novel method that takes compensating base changes across the binding sites into account. The algorithm works in two steps on two pre-generated multiple alignments. In the first step, individual base pairs with high reliability are found using the **PETfold** algorithm, which includes evolutionary and thermodynamic properties. In step two (where high reliability base pairs from step one are constrained as unpaired), the principle of cofolding is combined with hierarchical folding. The final prediction of *intra*- and *inter*-molecular base pairs consists of the reliabilities computed from the constrained expected accuracy scoring, which is an extended version of that used for individual multiple alignments.

**Results:** We derived a rather extensive algorithm. One of the advantages of our approach (in contrast to other RNA-RNA interaction prediction methods) is the application of covariance detection and prediction of pseudoknots between *intra*- and *inter*-molecular base pairs. As a proof of concept, we show an example and discuss the strengths and weaknesses of the approach.

## Background

Predicting RNA-RNA interactions is a rapidly growing area within RNA bioinformatics and is essential for the process of assigning function to known as well as *de novo* predicted non-coding RNAs (ncRNAs) such as those identified in *in silico* screens for RNA structures [1-7]. This candidate information along with the data generated from deep sequencing analyses emphasise the need to predict RNA-RNA interactions. In part, this is because there currently is no high-throughput method available for the reliable analysis of RNA-RNA interactions; however, computational prediction of RNA-RNA interactions is also essential for the identification of putative targets of

known and *de novo* predicted ncRNAs. With the main exception of microRNA target prediction, the current approaches essentially evaluate the stabilities of the common complexes between ncRNAs and target RNAs by computing the overall free energy using two major strategies (see, e.g., [8] for a recent review).

The first strategy, represented through the implementations of RNAup[9] and IntaRNA[10], uses pre-calculated values for all possible regions of interaction to determine the energy required to make that site accessible (called the ED-value for the energy difference). The ED-value is then used to calculate a combined energy of the energy given by the duplex formed by the two interaction regions and the ED-values of both interaction regions. RNAup has a complexity of  $O(n^3 + nw^5)$ , whereas IntaRNA has a complexity of  $O(n^2)$ , which makes it fast enough to be used in genome-wide screens. Both methods are able to predict complex interactions, like kissing

\* Correspondence: backofen@informatik.uni-freiburg.de

<sup>2</sup> Bioinformatics Group, University of Freiburg, Georges-Köhler-Allee 106, Freiburg, 79110, Germany

<sup>†</sup> Contributed equally

Full list of author information is available at the end of the article

hairpins, *as long as* the interaction is restricted to one region. However, there are well-known examples where several interaction sites were found, especially for longer ncRNAs. A prominent example is the interaction between *OxyS* and *fhlA* shown in [11].

The second strategy for RNA-RNA interaction predictions is usually handled with a class of approaches that simultaneously predict a common structure for both RNAs including their interaction. Some of the first approaches, *e.g.*, `pairfold`[12], `RNAcofold`[13] and the method presented by Dirks *et al.* as part of the `NUPack` package [14], concatenate the two sequences using a special linker character. Then, a modified version of the usual RNA folding methods (like `Mfold`[15] and `RNAfold`[16]) is applied to cope with the linker symbol to predict the correct energies. Otherwise, a loop containing the linker symbol would be treated like a hairpin or internal loop, leading to incorrect energy values.

The main disadvantage of the concatenation approach is that the set of candidate joint structures becomes restricted. For this reason, double kissing hairpin interactions (like in *OxyS-fhlA*) cannot be considered. However, alternative (but also most resource demanding) methods have been introduced and extend the class of allowed joint structures. The `IRIS` tool [17] allowed several kissing hairpins using a maximum number of base pair energy model. Then, Alkan *et al.* [18] presented a more realistic energy model and showed the NP-completeness of an unrestricted model. Both approaches predict structures with minimum free energy.

A more stable approach is to consider the partition function because it allows the calculation of interaction probabilities and melting temperatures. This problem was solved independently by Chitsaz *et al.* [19] and Huang *et al.* [20]. In [21], hybrid probabilities were calculated. These approaches have high time complexities of  $O(n^6)$ , which makes them infeasible for genome-wide applications. Methods to reduce the complexity range from approximation approaches [22,23] to sparsification of the dynamical programming matrix [24].

Here, we present an algorithm for the prediction of RNA-RNA interactions in existing multiple alignments of RNA sequences. Its rationale is based on the assumption that a non-negligible amount of the RNA-RNA interactions contain compensatory base changes across the binding sites. The algorithm presented herein is an extension of the `PETfold` algorithm [25] and makes further use of the principles from `RNAcofold` [13] and computational strategies for hierarchical folding, *e.g.* [26,27]. The latter approach was chosen due to the high computational costs of pseudoknot searches.

## Algorithm

The main idea of the introduced method is to use a hierarchical approach to predict an interaction by predicting reliable base pairs within a ncRNA and a mRNA (or another ncRNA), which is followed by prediction of reliable base pairs in the combined sequence. Via this approach, we are able to predict combined pseudoknotted structures, like kissing hairpins, that would be missed otherwise. In both steps, we apply a combined scoring method that predicts consensus base pairs from an alignment using evolutionary conservation and thermodynamic stability information.

The scoring for the first step is according to the standard `PETfold` approach, where we use thresholds for reliable base pairs that have been identified according to training on more than 30 verified interactions in bacteria, which is described later. For the second step, we define a constrained version of the `PETfold` scoring scheme.

Throughout this paper, we consider the concatenation of the two alignments and subsequently (in the base pair prediction process) the concatenation of the corresponding structures.  $\sigma$  will denote a set of base pairs, where the substructures in each part (*e.g.*, ncRNA, mRNA and the base pairs that participate in the interaction) in respective alignments are concatenated or nested (in the dot bracket notation, these substructures have alignment lengths of the ncRNA and mRNA respectively). We use  $(i, j)$  to denote a Watson-Crick or G-U wobble base pair between columns  $i$  and  $j$ . This base pair could be an *intra*-molecular pair in each of the RNA molecules (ncRNA or mRNA) or an *inter*-molecular pair that is involved in an interaction between molecules.

Depending on the context,  $\sigma$  will either be interpreted as a specific structure that implicitly defines the single-stranded positions or as a partial structure that describes an ensemble of structures. In the first case, we define the set of single-stranded positions of a sequence  $s$  as

$$ss(\sigma) = \left\{ i \mid \begin{array}{l} 1 \leq i \leq |s| \wedge \\ \forall j = 1, \dots, |s| : ((i, j) \notin \sigma \wedge (j, i) \notin \sigma) \end{array} \right\}.$$

In the second case, we use  $\mathcal{E}(\sigma) = \{\sigma' \mid \sigma' \supset \sigma\}$  to denote the ensemble of all specific structures  $\sigma'$  extending  $\sigma$ .  $\mathcal{S}(s)$  denotes the set of nested secondary structures that are defined for the sequence  $s$ . We use the same notation for the consensus structures of a given multiple alignment  $\mathcal{A}$  with  $n$  sequences  $s^1 \dots s^n$ . In this case, a position  $1 \leq i \leq |\mathcal{A}|$  refers to a column in the alignment. Furthermore, we use  $s \setminus \mathcal{A}$  to indicate a sequence  $s^1 \dots s^n$  from the alignment.

The algorithm, like **PETfold** is a maximum expected scoring approach that combines the evolutionary probabilities  $\Pr^{\text{ev}}[\sigma | \mathcal{A}]$  of a consensus structure,  $\sigma$ , given an alignment, with the thermodynamic probabilities of the associated structures in each sequence.  $\Pr^{\text{ev}}[\sigma | \mathcal{A}]$  is generated using the stochastic context-free grammar (SCFG) from the Pfold model [28]. The Pfold model allows the computation of the probability  $\Pr[\sigma | \mathcal{A}, T, M]$  of a consensus structure  $\sigma$  given an alignment  $\mathcal{A}$ , a phylogenetic tree  $T$  for that alignment, and a general background model  $M$  for secondary structures. Because the tree  $T$  is calculated from the alignment  $\mathcal{A}$ , and  $M$  is constant, we use  $\Pr^{\text{ev}}[\sigma | \mathcal{A}]$  as short for  $\Pr[\sigma | \mathcal{A}, T, M]$ .

The (secondary structure) model itself is based on a SCFG that provides a distribution of secondary structures for a given alignment. The combined probability of an alignment  $\mathcal{A}$  and a consensus structure  $\sigma$  is

$$\Pr[\mathcal{A}, \sigma | T, M] = \Pr[\mathcal{A} | T, \sigma] \Pr[\sigma | T, M],$$

where  $\Pr[\sigma | T, M]$  is the prior distribution of secondary structures and  $\Pr[\mathcal{A} | T, \sigma]$  is the probability of the alignment, given a known consensus structure. This is then transformed into  $\Pr[\mathcal{A}, \sigma | T, M]$  by applying the Bayesian rule, and further into the posterior distribution  $\Pr[\sigma | \mathcal{A}, T, M]$  of consensus structures  $\sigma$  by dividing by  $\Pr[\mathcal{A} | T, M]$ , which is the sum of all parse trees for an alignment  $\mathcal{A}$  given  $T$  and  $M$ . Note that the comma sign here is just a shortcut for  $\wedge$ , i.e.  $\Pr[A, B] = \Pr[A \wedge B]$ . We will still use  $\wedge$  where it is appropriate.

The probability distributions themselves are formed as follows. For  $\Pr[\mathcal{A} | T, \sigma]$ , there is an independent evaluation of all base pairs and single-stranded positions:

$$\Pr[\mathcal{A} | T, \sigma] = \prod_{(i,j) \in \sigma} \Pr_{\text{bp}}[\bar{\mathcal{A}}^i \bar{\mathcal{A}}^j | T] \times \prod_{i \in \text{ss}(\sigma)} \Pr_{\text{ss}}[\bar{\mathcal{A}}^i | T],$$

where  $\bar{\mathcal{A}}^i$  is the  $i$ th column of  $\mathcal{A}$ , and  $\sigma \notin \sigma_1^{\text{P}} \cup \sigma_2^{\text{P}}$  for the constrained folding, where  $\sigma_1^{\text{P}}$  ( $\sigma_2^{\text{P}}$  resp.) is the constrained structure on the first (second resp.) of the two concatenated alignments. For the prior model, the probability  $\Pr[\sigma | T, M]$  provides an overall distribution of the secondary structures, which is estimated from rRNA and tRNA sequences.  $M$  is given by the following simple SCFG:

$$S \rightarrow LS | L \quad F \rightarrow dFd | LS \quad L \rightarrow s | dFd.$$

The evolutionary model and the prior model for RNA structures used in the Pfold model are combined into a single SCFG that provides a distribution over  $\Pr[A, \sigma | T, M]$  (see additional file 1 for details).

To model the thermodynamic probabilities, we define  $\sigma(s^k, \mathcal{A})$  as the structure for the  $k$ -th sequence  $s^k$  of an alignment  $\mathcal{A}$  associated with the consensus structure  $\sigma$  of  $\mathcal{A}$ .  $\Pr^{\text{th}}[\sigma(s^k, \mathcal{A}) | s^k]$  is the corresponding thermodynamic probability as defined by McCaskill's partition function approach [29].

Using the maximum expected scoring approach, these probabilities are transformed into reliabilities in a two-step approach. Throughout the paper,  $\mathcal{R}_{\text{ss}}^{\ell}(i)$  is used to denote the reliability of a single-stranded region at alignment position  $i$  and  $\mathcal{R}_{\text{bp}}^{\ell}(i, j)$  the reliability of a consensus base pair  $(i, j)$ , where  $\ell = 1, 2$  refers to Step 1 or Step 2 of the combined approach.

#### Refresher: **PETfold** scoring

Here, we briefly recall the scoring of **PETfold**, which is a maximum expected accuracy scoring method. For simplicity, we will exclude a description of the scoring of single-stranded positions. However, they are scored the same way as in the original **PETfold** approach; for more details, see [25].

The **PETfold** score is the sum of the evolutionary accuracy values plus the average sum of the thermodynamic accuracy values. For the evolutionary part, we compute the expected accuracy (or overlap)  $\text{EA}^{\text{ev}}(\sigma)$  of a specific consensus structure  $\sigma$  with all possible consensus structures, which are weighted according to their probabilities:

$$\text{EA}^{\text{ev}}(\sigma) = \sum_{\sigma'} |\sigma \cap \sigma'| \times \Pr^{\text{ev}}[\sigma' | \mathcal{A}]. \quad (1)$$

Recall that  $\Pr^{\text{ev}}[\sigma' | \mathcal{A}]$  denotes the evolutionary probability of a structure  $\sigma'$  according to the Pfold SCFG as described above.  $|\sigma \cap \sigma'|$  is the number of base pairs that are common between  $\sigma$  and  $\sigma'$  and thus denotes the overlap between these two structures.

For the thermodynamic part, the expected accuracy  $\text{EA}^{\text{th}}(\sigma)$  of  $\sigma$  with all structures for all sequences according to the thermodynamic ensembles is defined by

$$\text{EA}^{\text{th}}(\sigma) = \sum_s \sum_{\sigma' \in S(s)} |\sigma \cap \sigma'| \times \Pr^{\text{th}}[\sigma'(s, \mathcal{A}) | s]. \quad (2)$$

The combined expected accuracy consists of both parts, generally weighted with 1 for the conservation portion and  $\beta$  for the thermodynamic accuracy:

$$EA(\sigma) = EA^{ev}(\sigma) + \frac{\beta}{n} \times EA^{th}(\sigma), \quad (3)$$

where  $n$  is the previously described number of sequences in the alignment. As shown previously [25], this final score can be calculated using the base pair reliabilities, where the combined reliability  $\mathcal{R}_{bp}(i, j)$  for a base pair  $(i, j)$  is given by

$$\begin{aligned} \mathcal{R}_{bp}(i, j) &= \sum_{\substack{\sigma' \text{ with} \\ (i, j) \in \sigma'}} 1 \times \Pr^{ev}[\sigma' | \mathcal{A}] \\ &+ \frac{\beta}{n} \times \sum_s \sum_{\substack{\sigma' \text{ with} \\ (i, j) \in \sigma'}} 1 \times \Pr^{th}[\sigma'(s, \mathcal{A}) | s] \quad (4) \\ &= \mathcal{R}_{bp}^{ev}(i, j, \mathcal{A}) + \frac{\beta}{n} \times \sum_s \Pr_{bp}^{th}(i, j, s), \end{aligned}$$

where  $\Pr_{bp}^{th}(i, j, s)$  is the base pair probability of the pair  $(k, l)$  associated with columns  $(i, j)$  in sequence  $s$ . These reliabilities are calculated with an inside/outside algorithm and are central to the hierarchical approach presented in the following sections. The expected accuracy can then be calculated from the base pair reliabilities by

$$EA(\sigma) = \sum_{(i, j) \in \sigma} \mathcal{R}_{bp}(i, j). \quad (5)$$

The consensus structure with the maximal reliability is then calculated using a Nussinov-style algorithm [30], where the base pairs are evaluated with reliabilities.

### Step 1: Intra-molecular partial structures

We use two alignments  $\mathcal{A}_1$  and  $\mathcal{A}_2$  of sequences  $s_1^1 \dots s_1^n$  and  $s_2^1 \dots s_2^m$ , where  $s_1^k$  is a ncRNA and  $s_2^k$  is its target sequence. For convenience, we adopt the convention of `RNAcofold` and assume that the positions in  $s_1^k$  are numbered  $1 \leq i \leq |s_1^k|$  and the positions in  $s_2^k$  are numbered  $|s_1^k| + 1 \leq i \leq |s_1^k| + |s_2^k|$ .

#### Selection of the initial structure

In the first step of the pipeline, we obtain the base pair reliabilities from Equation (4), which we denote  $\mathcal{R}_{bp}^1(i, j)$ . Using these reliabilities, the partial (constrained) struc-

tures  $\sigma_1^p$  and  $\sigma_2^p$  are determined independently for  $\mathcal{A}_1$  and  $\mathcal{A}_2$ . In the following steps, let  $\mathcal{A}$  be either  $\mathcal{A}_1$  or  $\mathcal{A}_2$  and  $\sigma^p$  be the partial structure calculated for  $\mathcal{A}$ . This is done by selecting only base pairs  $(i, j)$  with

$$\mathcal{R}_{bp}^1(i, j) \geq \delta,$$

where  $\delta$  is a cut-off that must be  $\geq 0.5$  to avoid crossing structures. This is similar to the method by which consensus structures are predicted for single sequences [31] and has been shown to be more reliable for the prediction of consensus structures from alignments.

Here, however, we also have to estimate the contribution of each of the partial structures to the complete solution. Because the set of base pairs from a predicted consensus structure do not necessarily form a reasonable structure, we account for this by introducing a second threshold  $\gamma$ . High values for this threshold guarantee that each sequence used to create the consensus structure has a high likelihood and that the approximation, which we apply in the second step (as will be described by Equation (14)), is accurate.

To find the optimal value of the reliability threshold  $\delta$ , its value is increased until the resulting ensemble of structures  $\mathcal{E}(\sigma^p)$  that are compatible with the partial structure  $\sigma^p$  is probable enough in the evolutionary model, in the thermodynamic model, or in both models, which is when

$$\Pr^{ev}[\mathcal{E}(\sigma^p) | \mathcal{A}] \geq \gamma \text{ or } \frac{1}{n} \sum_{s \in \mathcal{A}} \Pr^{th}[\mathcal{E}(\sigma^p(s, \mathcal{A})) | s] \geq \gamma.$$

Here,  $\Pr^{ev}[\mathcal{E}(\sigma^p) | \mathcal{A}]$  ( $= \Pr^{ev}[\mathcal{E}(\sigma^p) | \mathcal{A}, T, M]$ ) is the probability of the partial structure  $\sigma^p$  given the alignment  $\mathcal{A}$  in the evolutionary model  $M$  and tree  $T$ . This can be calculated from `Pfold` with the SCFG that combines the prior structural model with evolutionary information from the alignment (see additional file 1) as follows:

$$\begin{aligned} &\Pr^{ev}[\mathcal{E}(\sigma^p) | \mathcal{A}, T, M] \\ &= \frac{\Pr[\mathcal{E}(\sigma^p) | \mathcal{A}, T, M] \times \Pr[\mathcal{A} | T, M]}{1 \times \Pr[\mathcal{A} | T, M]} \quad (6) \\ &= \frac{\Pr[\mathcal{E}(\sigma^p), \mathcal{A} | T, M]}{\Pr[\mathcal{A} | T, M]} \end{aligned}$$

The term  $\Pr[\mathcal{A} | T, M]$  has already been calculated (personal communication with Bjarne Knudsen) in `Pfold` as

the sum of all possible parse trees for an alignment  $\mathcal{A}$ , given  $T, M$ :

$$\Pr[\mathcal{A} \mid T, M] = \sum_{\sigma} \Pr[\mathcal{A}, \sigma \mid T, M].$$

Here, we add the calculation of

$$\Pr[\mathcal{E}(\sigma^P), \mathcal{A} \mid T, M] = \sum_{\sigma \in \mathcal{E}(\sigma^P)} \Pr[\mathcal{A}, \sigma \mid T, M]$$

to `Pfold` by summing over all possible parse trees that are compatible with  $\sigma$ .

$\Pr^{\text{th}}[\mathcal{E}(\sigma^P(s, \mathcal{A})) \mid s]$  is the probability of the partial structure  $\sigma^P$  given a sequence  $s$  in the thermodynamic model. This probability can be calculated using constrained partition folding as follows:

$$\begin{aligned} \Pr^{\text{th}}[\mathcal{E}(\sigma^P(s, \mathcal{A})) \mid s] &= \frac{e^{-E^S} \mathcal{E}(\sigma^P(s, \mathcal{A}))}{e^{-E^S_{\text{all}}} RT} \\ &= e^{\frac{E^S_{\text{all}} - E^S}{RT}} \frac{\mathcal{E}(\sigma^P(s, \mathcal{A}))}{RT}, \end{aligned} \quad (7)$$

where  $E^S_{\text{all}}$  is the free energy of the whole ensemble (as determined by `RNAfold` with parameters `-p -d2`) and  $E^S_{\mathcal{E}(\sigma^P(s, \mathcal{A}))}$  is the free energy of the ensemble of structures  $\mathcal{E}(\sigma^P(s, \mathcal{A}))$  with the base pairs in  $\sigma^P(s, \mathcal{A})$  as constraints, which can be calculated by `RNAfold` with parameters `-C -p -d2`.

#### Extension of constrained stems

Reliable *intra*-molecular base pairs are constrained as single-stranded in Step 2 of the algorithm because we are interested in pseudoknots of the concatenated sequence and the interactions in these induced loop regions. The drawback of this *ansatz* is that *intra*-molecular stems get instable because of intermediate unbased constraints. Thus, we may get incomplete stems. To deal with this problem, we extend the constrained stems. Inner and outer base pairs are added as long as the average reliability of the inner or outer extended stem, respectively, is larger than the threshold  $\delta$ , and the probability of the partial structure is greater than or equal to  $\gamma$  either in the evolutionary or the thermodynamic model. That is, the average reliability of the total, extended stem has to be larger than a threshold.

Step 1 is summarised as pseudocode in Figure 1.

```

for Alignment  $\mathcal{A}_1, \mathcal{A}_2$  do
  calculate tree  $T$ ,
  phylogenetic reliabilities  $\mathcal{R}^{1,\text{ev}}$ ,
  thermodynamic probabilities  $\Pr^{1,\text{th}}$ 
   $\implies \mathcal{R}_{\text{bp}}^1(i, j), \mathcal{R}_{\text{ss}}^1(i) \leftarrow$  PETfold model
  repeat
    for all  $(i, j)$  do
      if  $\mathcal{R}_{\text{bp}}^1(i, j) \geq \delta$  then
        add base pair  $(i, j)$  to  $\sigma^P$ 
      end if
    end for
    calculate partial structure probabilities
       $\Pr^{\text{ev}}[\mathcal{E}(\sigma^P) \mid \mathcal{A}]$  and  $\Pr^{\text{th}}[\mathcal{E}(\sigma^P) \mid s]$ 
    increase  $\delta$ 
  until
     $\Pr^{\text{ev}}[\mathcal{E}(\sigma^P) \mid \mathcal{A}] \geq \gamma \parallel \frac{1}{n} \sum_s \Pr^{\text{th}}[\mathcal{E}(\sigma^P) \mid s] \geq \gamma$ 

  for all stem  $\mathcal{S} \subset \sigma^P$  do
    for  $adjacent = (\text{inner}, \text{outer})$  do
      repeat
         $b = adjacent$  base pair of  $\mathcal{S}$ 
         $\mathcal{S}_{\text{old}} = \mathcal{S}; \mathcal{S} = \mathcal{S} \cup \{b\}$ 
         $\sigma_{\text{old}}^P = \sigma^P; \sigma^P = \sigma^P \cup \mathcal{S}$ 
        calculate  $\Pr^{\text{ev}}[\mathcal{E}(\sigma^P) \mid \mathcal{A}], \Pr^{\text{th}}[\mathcal{E}(\sigma^P) \mid s]$ 
      until average  $\mathcal{R}_{\text{bp}}^1$  of  $\mathcal{S} < \delta \parallel$ 
       $(\Pr^{\text{ev}}[\mathcal{E}(\sigma^P) \mid \mathcal{A}] < \gamma \ \&\& \ \frac{1}{n} \sum_s \Pr^{\text{th}}[\mathcal{E}(\sigma^P) \mid s] < \gamma)$ 
       $\mathcal{S} = \mathcal{S}_{\text{old}}; \sigma^P = \sigma_{\text{old}}^P$ 
    end for
  end for
end for
    
```

**Figure 1** Pseudocode for Step 1.

#### Step 2: Constrained expected accuracy scoring

In the following,  $s_1 \& s_2$  denote the concatenated sequences of the two sequences  $s_1, s_2$  using the additional linker symbol  $\&$  as done in `RNAcofold`. For Step 2 of the scoring, we calculate the expected accuracy of the ensemble of structures  $\sigma$  of  $s_1 \& s_2$ , which constitutes an interaction under the constraint that  $\sigma$  contains the partial reliable structures  $\sigma_1^P$  and  $\sigma_2^P$  of  $s_1$  and  $s_2$ , respectively. Because we use the numbering convention of `RNAcofold`, the union  $\sigma_1^P \cup \sigma_2^P$  of the two partial structures  $\sigma_1^P$  and  $\sigma_2^P$  is the partial structure of  $s_1 \& s_2$ .

Now we have two problems to solve. On the one hand, we want to calculate the constrained accuracy given the partial structures  $\sigma_1^P$  and  $\sigma_2^P$ , which is defined as

$$EA_{\sigma_1^P, \sigma_2^P}(\sigma) = EA_{\sigma_1^P, \sigma_2^P}^{\text{ev}}(\sigma) + \frac{\beta}{n} \times EA_{\sigma_1^P, \sigma_2^P}^{\text{th}}(\sigma). \quad (8)$$

On the other hand, we have to find a combined score for the partial structures  $\sigma_1^p$  and  $\sigma_2^p$ , and the interaction  $\sigma_{\text{int}}$  to evaluate the quality of an predicted interaction. The score must be maximal according to Equation (8).

We will demonstrate the problem and our solution for the thermodynamic folding. However, the same analysis applies to the evolutionary part, which is described later.

**The thermodynamic part**

The simplest formal solution to this problem would be to investigate directly the expected accuracy of joint structures  $\sigma$ :

$$\begin{aligned} EA^{\text{th}}(\sigma) &= \sum_{s_1 \& s_2} \sum_{\sigma'} |\sigma(s_1 \& s_2, \mathcal{A}) \cap \sigma'| \times \Pr^{\text{th}}[\sigma' | s_1 \& s_2] \\ &= \sum_{s_1 \& s_2 \in \mathcal{A}} EA_{s_1 \& s_2}^{\text{th}}(\sigma(s_1 \& s_2, \mathcal{A})), \end{aligned}$$

where  $EA_{s_1 \& s_2}^{\text{th}}(\sigma)$  is the expected accuracy of a structure in one sequence pair  $s_1 \& s_2 \in \mathcal{A}$ .

However, this would require that we compute the distribution  $\Pr^{\text{th}}[\sigma | s_1 \& s_2]$ , which can be done by a partition function approach for interacting structures. This is NP-complete in the full model [18] and even  $O(n^6)$  in a restricted model [19,20], which is why the two-step approach is necessary. In the following, we ignore the index "th" for simplicity.

The relationship between  $EA_{\sigma_1^p, \sigma_2^p}(\sigma)$  and  $EA(\sigma)$  is now quantified. In the following, for a structure  $\sigma$ , we use  $\sigma_1 \cup \sigma_2 \cup \sigma_{\text{int}}$  to denote the partition of the base pairs of the first sequence,  $\sigma_1$ , the base pairs of the second sequence,  $\sigma_2$ , and the interacting base pairs,  $\sigma_{\text{int}}$ . Furthermore, for the partial structure  $\sigma$ , we use  $\mathcal{E}_1(\sigma)$  to denote the set of structures that extends  $\sigma$  using base pairs within the first sequence, *i.e.*,

$$\begin{aligned} \mathcal{E}_1(\sigma) &= \{\sigma' \supseteq \sigma \mid \forall (i, j) \in \sigma' \setminus \sigma : 1 \leq i < j \leq |s_1|\} \\ \mathcal{E}_2(\sigma) &= \{\sigma' \supseteq \sigma \mid \forall (i, j) \in \sigma' \setminus \sigma : \\ &\quad |s_1| + 1 \leq i < j \leq |s_2|\} \\ \mathcal{E}_{\text{int}}(\sigma) &= \{\sigma' \supseteq \sigma \mid \forall (i, j) \in \sigma' \setminus \sigma : \\ &\quad 1 \leq i \leq |s_1| \wedge |s_1| + 1 \leq j \leq |s_2|\} \end{aligned}$$

The ensembles  $\mathcal{E}_{1,\text{int}}(\sigma)$ ,  $\mathcal{E}_{2,\text{int}}(\sigma)$  and  $\mathcal{E}_{1,2}(\sigma)$  are defined analogously.

Our approach uses one simplification, namely the assumption that the reliabilities for *intra*-molecular base pairs are dominated by the *intra*-molecular folding. This

is equivalent to the assumption that the two structures fold independently. We formulate this as follows:

$$\Pr[\sigma_1 \mid \sigma_2, s_1 \& s_2] = \Pr[\sigma_1 \mid s_1].$$

Because  $\sigma_1$  and  $\sigma_2$  are *partial joint* structures, this can be written using the ensemble function

$$\Pr[\mathcal{E}_{2,\text{int}}(\sigma_1) \mid \mathcal{E}_{1,\text{int}}(\sigma_2), s_1 \& s_2] = \Pr[\sigma_1 \mid s_1]. \quad (9)$$

The implication of this assumption is that the probabilities of the two structures  $\sigma_1$  and  $\sigma_2$  are merged independently into the joint probability  $\Pr[\mathcal{E}_{\text{int}}(\sigma_1 \cup \sigma_2) | s_1 \& s_2]$ , see Equation (11) below. First, note that for two partial structures

$$\Pr[\mathcal{E}(\sigma^p \cup \sigma^{p'}) \mid s] = \Pr[\mathcal{E}(\sigma^p) \wedge \mathcal{E}(\sigma^{p'}) \mid s],$$

by definition. Hence,

$$\begin{aligned} \Pr[\mathcal{E}_{\text{int}}(\sigma_1 \cup \sigma_2) \mid s_1 \& s_2] &= \Pr[\mathcal{E}_{2,\text{int}}(\sigma_1) \wedge \mathcal{E}_{1,\text{int}}(\sigma_2) \mid s_1 \& s_2] \\ &= \Pr[\mathcal{E}_{2,\text{int}}(\sigma_1) \mid \mathcal{E}_{1,\text{int}}(\sigma_2), s_1 \& s_2] \\ &\quad \times \Pr[\mathcal{E}_{1,\text{int}}(\sigma_2) \mid s_1 \& s_2] \\ &\stackrel{Eq.(9)}{=} \Pr[\sigma_1 \mid s_1] \times \Pr[\mathcal{E}_{1,\text{int}}(\sigma_2) \mid s_1 \& s_2]. \end{aligned}$$

Intuitively,  $\Pr[\mathcal{E}_{1,\text{int}}(\sigma_2) | s_1 \& s_2]$  should be the same as  $\Pr[\sigma_2 | s_2]$ . This can be derived using the total probability formula:

$$\begin{aligned} \Pr[\mathcal{E}_{1,\text{int}}(\sigma_2) \mid s_1 \& s_2] &= \sum_{\sigma_1} \left( \Pr[\mathcal{E}_{1,\text{int}}(\sigma_2) \mid \mathcal{E}_{2,\text{int}}(\sigma_1), s_1 \& s_2] \right. \\ &\quad \left. \times \Pr[\mathcal{E}_{2,\text{int}}(\sigma_1) \mid s_1 \& s_2] \right) \\ &\stackrel{Eq.(9)}{=} \sum_{\sigma_1} \Pr[\sigma_2 \mid s_2] \times \Pr[\mathcal{E}_{2,\text{int}}(\sigma_1) \mid s_1 \& s_2] \\ &= \Pr[\sigma_2 \mid s_2] \times \sum_{\sigma_1} \Pr[\mathcal{E}_{2,\text{int}}(\sigma_1) \mid s_1 \& s_2] \\ &= \Pr[\sigma_2 \mid s_2] \times 1 \end{aligned} \quad (10)$$

Combining these equations we obtain the independence property:

$$\Pr[\mathcal{E}_{\text{int}}(\sigma_1 \cup \sigma_2) \mid s_1 \& s_2] = \Pr[\sigma_1 \mid s_1] \times \Pr[\sigma_2 \mid s_2] \quad (11)$$

Now we will use this property to relate  $EA_{\sigma_1^p, \sigma_2^p}(\sigma)$  to  $EA(\sigma)$ . The independence property, as described in Equation (9), and the additivity of the expectation is the implication of the expected accuracy of a joint structure, which is the sum of the expected accuracy of the *intra*-molecular structures and the expected accuracy of the *inter*-molecular portion. To illustrate this, note that for any  $\sigma, \sigma'$

$$|\sigma \cap \sigma'| = |\sigma_1 \cap \sigma'_1| + |\sigma_2 \cap \sigma'_2| + |\sigma_{\text{int}} \cap \sigma'_{\text{int}}|$$

by definition. Hence, by the additivity of the expectation we get

$$\begin{aligned} EA_{s_1 \& s_2}^{\text{th}}(\sigma) &= \sum_{\sigma'} |\sigma \cap \sigma'| \times \Pr^{\text{th}}[\sigma' | s_1 \& s_2] \\ &= \sum_{\sigma'} |\sigma_1 \cap \sigma'_1| \times \Pr[\sigma' | s_1 \& s_2] \\ &\quad + \sum_{\sigma'} |\sigma_2 \cap \sigma'_2| \times \Pr[\sigma' | s_1 \& s_2] \\ &\quad + \sum_{\sigma'} |\sigma_{\text{int}} \cap \sigma'_{\text{int}}| \times \Pr[\sigma' | s_1 \& s_2]. \end{aligned}$$

Now we can rewrite the first term  $\sum_{\sigma'} |\sigma_1 \cap \sigma'_1| \times \Pr[\sigma' | s_1 \& s_2]$  using the independence property as follows:

$$\begin{aligned} &\sum_{\sigma'} |\sigma_1 \cap \sigma'_1| \times \Pr[\sigma' | s_1 \& s_2] \\ &= \sum_{\sigma'_1} \sum_{\sigma'_2, \sigma'_{\text{int}}} |\sigma_1 \cap \sigma'_1| \times \Pr[\sigma'_1 \cup \sigma'_2 \cup \sigma'_{\text{int}} | s_1 \& s_2] \\ &= \sum_{\sigma'_1} |\sigma_1 \cap \sigma'_1| \times \Pr[\mathcal{E}_{2, \text{int}}(\sigma'_1) | s_1 \& s_2] \\ &\stackrel{\text{Eq. (10)}}{=} \sum_{\sigma'_1} |\sigma_1 \cap \sigma'_1| \times \Pr[\sigma'_1 | s_1], \end{aligned}$$

which is the expected accuracy of  $\sigma_1$  in the sequence  $s_1$ . Analogously, we can do this for the second term  $\sum_{\sigma'} |\sigma_2 \cap \sigma'_2| \times \Pr[\sigma' | s_1 \& s_2]$ . Thus,  $EA_{s_1 \& s_2}^{\text{th}}(\sigma)$  is the sum of the expected accuracies in the first and the second sequences and the expected accuracy of the interaction:

$$\begin{aligned} EA_{s_1 \& s_2}^{\text{th}}(\sigma) &= EA_{s_1}^{\text{th}}(\sigma_1) + EA_{s_2}^{\text{th}}(\sigma_2) \\ &\quad + \sum_{\sigma'} |\sigma_{\text{int}} \cap \sigma'_{\text{int}}| \times \Pr[\sigma' | s_1 \& s_2]. \end{aligned} \tag{12}$$

For the expected accuracy of the interaction

$$EA^{\text{th, int}}(\sigma) = \sum_{\sigma'} |\sigma_{\text{int}} \cap \sigma'_{\text{int}}| \times \Pr[\sigma' | s_1 \& s_2] \tag{13}$$

we still need to define  $\Pr[\sigma | s_1 \& s_2]$ . For every  $\sigma = \sigma_1 \cup \sigma_2 \cup \sigma_{\text{int}}$

$$\begin{aligned} &\Pr[\sigma_1 \cup \sigma_2 \cup \sigma_{\text{int}} | s_1 \& s_2] \\ &= \Pr[\mathcal{E}_{\text{int}}(\sigma_1 \cup \sigma_2) \wedge \mathcal{E}_{1,2}(\sigma_{\text{int}}) | s_1 \& s_2] \\ &= \Pr[\mathcal{E}_{1,2}(\sigma_{\text{int}}) | \mathcal{E}_{\text{int}}(\sigma_1 \cup \sigma_2), s_1 \& s_2] \\ &\quad \times \Pr[\mathcal{E}_{\text{int}}(\sigma_1 \cup \sigma_2) | s_1 \times s_2] \\ &\stackrel{\text{Eq. (11)}}{=} \Pr[\mathcal{E}_{1,2}(\sigma_{\text{int}}) | \mathcal{E}_{\text{int}}(\sigma_1 \cup \sigma_2), s_1 \& s_2] \\ &\quad \times \Pr[\sigma_1 | s_1] \times \Pr[\sigma_2 | s_2] \end{aligned}$$

Thus, in principle, to calculate the expected accuracy  $EA^{\text{th, int}}(\sigma)$  for the interaction, we must sum over all structures in  $\sigma_1$  and  $\sigma_2$ :

$$\begin{aligned} EA^{\text{th, int}}(\sigma) &= \sum_{\sigma_{\text{int}}} \sum_{\sigma'_1, \sigma'_2} |\sigma_{\text{int}} \cap \sigma'_{\text{int}}| \times \Pr[\sigma'_{\text{int}} \cup \sigma'_1 \cup \sigma'_2 | s_1 \& s_2] \\ &= \sum_{\sigma_{\text{int}}} |\sigma_{\text{int}} \cap \sigma'_{\text{int}}| \times \sum_{\sigma'_1, \sigma'_2} \Pr[\sigma'_{\text{int}} \cup \sigma'_1 \cup \sigma'_2 | s_1 \& s_2] \end{aligned}$$

Because this is not feasible, we restrict ourselves to an ensemble of structures. Thus, instead of summing over all possible  $\sigma_1$  and  $\sigma_2$ , we use the partial structures  $\sigma_1^p$  and  $\sigma_2^p$  that were determined in the first step and approximate  $EA^{\text{th, int}}(\sigma)$  by

$$\begin{aligned} EA^{\text{th, int}}(\sigma) &= \sum_{\sigma_{\text{int}}} |\sigma_{\text{int}} \cap \sigma'_{\text{int}}| \times \sum_{\substack{\sigma'_1 \in \mathcal{E}(\sigma_1^p) \\ \sigma'_2 \in \mathcal{E}(\sigma_2^p)}} \Pr[\sigma'_{\text{int}} \cup \sigma'_1 \cup \sigma'_2 | s_1 \& s_2] \end{aligned}$$

The second sum can now be simplified as follows:

$$\begin{aligned}
 & \sum_{\substack{\sigma'_1 \in \mathcal{E}(\sigma_1^p) \\ \sigma'_2 \in \mathcal{E}(\sigma_2^p)}} \Pr[\sigma'_{\text{int}} \cup \sigma'_1 \cup \sigma'_2 \mid s_1 \& s_2] \\
 &= \Pr[\mathcal{E}_{1,2}(\sigma'_{\text{int}} \cup \sigma_1^p \cup \sigma_2^p) \mid s_1 \& s_2] \\
 &= \Pr[\mathcal{E}_{1,2}(\sigma'_{\text{int}}) \wedge \mathcal{E}(\sigma_1^p \cup \sigma_2^p) \mid s_1 \& s_2] \\
 &= \Pr[\mathcal{E}_{1,2}(\sigma'_{\text{int}}) \mid \mathcal{E}(\sigma_1^p \cup \sigma_2^p), s_1 \& s_2] \\
 &\quad \times \Pr[\mathcal{E}(\sigma_1^p \cup \sigma_2^p) \mid s_1 \& s_2] \\
 &\stackrel{\text{Eq. (11')}}{=} \Pr[\mathcal{E}_{1,2}(\sigma'_{\text{int}}) \mid \mathcal{E}(\sigma_1^p \cup \sigma_2^p), s_1 \& s_2] \\
 &\quad \times \Pr[\mathcal{E}_1(\sigma_1^p) \mid s_1] \times \Pr[\mathcal{E}_2(\sigma_2^p) \mid s_2],
 \end{aligned}$$

where Equation (11') indicates the variation of the independence assumption of Equation (11) for the structure ensembles (see additional file 1). Thus, we finally have

$$\begin{aligned}
 & \text{EA}^{\text{th,int}'}(\sigma) \\
 &= \sum_{\sigma'_{\text{int}}} \left( \Pr[\sigma_{\text{int}} \cap \sigma'_{\text{int}} \mid \Pr[\mathcal{E}_1(\sigma_1^p) \mid s_1] \times \Pr[\mathcal{E}_2(\sigma_2^p) \mid s_2]] \right. \\
 &\quad \left. \times \Pr[\mathcal{E}_{1,2}(\sigma'_{\text{int}}) \mid \mathcal{E}(\sigma_1^p \cup \sigma_2^p), s_1 \& s_2] \right) \quad (14)
 \end{aligned}$$

Now  $\Pr[\mathcal{E}_{1,2}(\sigma'_{\text{int}}) \mid \mathcal{E}(\sigma_1^p \cup \sigma_2^p), s_1 \& s_2]$  is the constrained folding, where the positions covered by  $\sigma_1^p$  and  $\sigma_2^p$  are fixed. However, we have the problem that these structures might contain pseudoknots. Recall that the positions in  $\sigma_1^p$  and  $\sigma_2^p$  are fixed for folding and that we are considering all structures  $\sigma$  that contain  $\sigma_1^p \cup \sigma_2^p$  and are nested on  $\sigma_{\text{int}} = \sigma \setminus (\sigma_1^p \cup \sigma_2^p)$ . Technically, we solve the problem using the fact that the set of structures that is nested on  $\sigma_{\text{int}}$  and compatible with  $\sigma_1^p \cup \sigma_2^p$  is selected by considering all structures where the positions of  $\sigma_1^p \cup \sigma_2^p$  are constrained as single-stranded. This implies that we use constrained cofolding via RNACoFold (parameters -C -p -d2), and the constraint ...  $x_1 x_1 \dots$  & ...  $x_2 x_2 \dots$ , where  $x_1$  (resp.  $x_2$ ) denotes a position from  $\sigma_1^p$  (resp.  $\sigma_2^p$ ) that is constrained as single-stranded. The main difference is that the energy contributions could be slightly different, and therefore, we obtain only an approximation of the real distribution. For example, an extension of a helix in  $\sigma_1^p$  would be evaluated as an internal loop or hairpin. Note that this is not a major problem because we are mainly interested in the

*inter*-molecular base pairs between  $s_1$  and  $s_2$  in this step. However, the recursion scheme of RNACoFold could easily be adapted to use new symbols for base pair constraints and a scoring scheme that is common to hierarchical approaches of pseudoknot structure prediction, which would avoid these problems.

Finally, we can rewrite the thermodynamic accuracy as the sum of probabilities as indicated in Equation (5). As shown in Equation (12), for a base pair  $(i, j) \in \sigma_\ell^p$  ( $\ell = 1, 2$ ), we want to use the probability of the associated sequence. To avoid competition with the probabilities for the *intra*-molecular base pairs calculated from RNACoFold, we set all of these base pairs to the same probability  $\Pr^{\text{th}}[\mathcal{E}_\ell(\sigma_\ell^p(s_\ell, \mathcal{A})) \mid s_\ell]$  as described in Equation (7). For the *inter*-molecular base pairs, we use the base pair probabilities as provided by RNACoFold with constraints, which model  $\Pr[\sigma' \mid \mathcal{E}(\sigma_1^p \cup \sigma_2^p) \wedge s_1 \& s_2]$  from the constrained cofolding. However, these raw base pair probabilities (in the following denoted by  $\Pr_{\text{bp,raw}}^{2,\text{th}}(i, j)$ ) are calculated under the constraint of  $\sigma_1^p \cup \sigma_2^p$  and have therefore (to obtain the final base pair probabilities) to be multiplied by  $\Pr[\mathcal{E}_1(\sigma_1^p) \mid s_1] \times \Pr[\mathcal{E}_2(\sigma_2^p) \mid s_2]$  as indicated by Equation (14). Thus, we can score each base pair as follows:

$$\begin{aligned}
 & \Pr_{\text{bp}}^{2,\text{th}}(i, j, s_1 \& s_2) \\
 &= \begin{cases} 1 \times \Pr_{\text{bp}}^{\text{th}}(i, j, s_1) & \text{if } 1 \leq i, j \leq s_1 \mid \\ 1 \times \Pr_{\text{bp}}^{\text{th}}(i, j, s_2) & \text{if } |s_1| + 1 \leq i, j \leq s_2 \mid \\ \Pr_{\text{bp,raw}}^{2,\text{th}}(i, j) \times \prod_{\ell=1,2} \Pr^{\text{th}}[\mathcal{E}_\ell(\sigma_\ell^p(s_\ell, \mathcal{A})) \mid s_\ell] & \end{cases} \quad (15)
 \end{aligned}$$

where the 1 reflects the fixed reliability. However, we deviate from this scoring to weaken the independence assumption for the *intra*-molecular base pairs, which allows us to determine new *intra*-molecular base pairs from the constrained application of RNA-cofold. Thus, we score only the base pairs from the partial structures  $\sigma_1^p$  and  $\sigma_2^p$  with the probability in the associated sequence. In addition, to avoid competition with the probabilities for these base pairs calculated from RNACoFold, we simply set all of these base pairs to the same probability  $\Pr^{\text{th}}[\mathcal{E}_\ell(\sigma_\ell^p(s_\ell, \mathcal{A})) \mid s_\ell]$  as described in Equation (7). To summarise, given the partial consensus structures  $\sigma_1^p$  and  $\sigma_2^p$  for an alignment  $\mathcal{A}_1$  &  $\mathcal{A}_2$  as cal-



culated in Step 1, the probability for a base pair  $(i, j)$  in sequence  $s_1 \& s_2 \in \mathcal{A}_1 \& \mathcal{A}_2$  in the second step is:

$$\begin{aligned} & \Pr_{\text{bp}}^{2,\text{th}}(i, j, s_1 \& s_2) \\ &= \begin{cases} 1 \times \Pr^{\text{th}}[\mathcal{E}_1(\sigma_1^{\text{p}}(s_1, \mathcal{A})) | s_1] & \text{if } (i, j) \in \sigma_1^{\text{p}} \\ 1 \times \Pr^{\text{th}}[\mathcal{E}_2(\sigma_2^{\text{p}}(s_2, \mathcal{A})) | s_2] & \text{if } (i, j) \in \sigma_2^{\text{p}} \\ \Pr_{\text{bp,raw}}^{2,\text{th}}(i, j) \times \prod_{\ell=1,2} \Pr^{\text{th}}[\mathcal{E}_\ell(\sigma_\ell^{\text{p}}(s_\ell, \mathcal{A})) | s_\ell] & \text{else} \end{cases} \end{aligned} \quad (16)$$

### Single-stranded probabilities

Single-stranded probabilities are integrated in a similar way as the base pair probabilities, but with different weighting. The single-stranded probabilities are as follows:

$$\begin{aligned} & \Pr_{\text{ss}}^{2,\text{th}}(i, s_1 \& s_2) \\ &= \begin{cases} 0 & \text{if } \exists j \text{ with } (i, j) \in \sigma_1^{\text{p}} \text{ or } (j, i) \in \sigma_1^{\text{p}} \\ 0 & \text{if } \exists j \text{ with } (i, j) \in \sigma_2^{\text{p}} \text{ or } (j, i) \in \sigma_2^{\text{p}} \\ \Pr_{\text{ss,raw}}^{2,\text{th}}(i) \times \prod_{\ell=1,2} \Pr^{\text{th}}[\mathcal{E}_\ell(\sigma_\ell^{\text{p}}(s_\ell, \mathcal{A})) | s_\ell] & \text{else} \end{cases} \end{aligned} \quad (17)$$

Given the structure  $\sigma$  on an alignment  $\mathcal{A}$  with  $m$  columns, the set of all single-stranded positions in the consensus structure is denoted as  $\text{ss}(\sigma) = \{i \notin \sigma | 1 \leq i \leq m\}$ . Taking this into consideration, the complete version of Equation (2) is

$$\begin{aligned} & \text{EA}^{\text{th}}(\sigma) \\ &= \sum_{s, \sigma'} [|\sigma \cap \sigma'| + \alpha | \text{ss}(\sigma) \cap \text{ss}(\sigma') |] \times \Pr^{\text{th}}[\sigma'(s, \mathcal{A}) | s] \end{aligned}$$

and the evolutionary accuracy is determined similarly. The combined score is the sum of the base pair reliabilities and single-stranded reliabilities (weighted with the parameter  $\alpha$ ). For details, see [25].

### The evolutionary part

The calculation for the presented thermodynamic accuracy is purely based on constrained folding. To obtain the complete constrained folding, we use the same approach for the evolutionary accuracy by applying a version of Pfold[28] that incorporates the constraints. For that purpose, the raw structural reliabilities  $\mathcal{R}_{\text{bp,raw}}^{2,\text{ev}}(i, j)$  and

$\mathcal{R}_{\text{ss,raw}}^{2,\text{ev}}(i)$  are calculated by the constrained folding with Pfold using the phylogenetic tree deduced from the concatenated alignment. As a linker, three prior-free columns are inserted between both alignments. The evolu-

tionary reliabilities  $\mathcal{R}_{\text{bp}}^{2,\text{ev}}(i, j)$  for a base pair  $(i, j)$  and  $\mathcal{R}_{\text{ss}}^{2,\text{ev}}(i)$  for a single-stranded position  $i$  are calculated in the same manner as  $\Pr_{\text{bp}}^{2,\text{th}}(i, j, s_1 \& s_2)$  in Equation (16):

$$\begin{aligned} & \mathcal{R}_{\text{bp}}^{2,\text{ev}}(i, j, \mathcal{A}_1 \& \mathcal{A}_2) \\ &= \begin{cases} 1 \times \Pr^{\text{ev}}[\mathcal{E}_1(\sigma_1^{\text{p}}) | \mathcal{A}_1] & \text{if } (i, j) \in \sigma_1^{\text{p}} \\ 1 \times \Pr^{\text{ev}}[\mathcal{E}_2(\sigma_2^{\text{p}}) | \mathcal{A}_2] & \text{if } (i, j) \in \sigma_2^{\text{p}} \\ \mathcal{R}_{\text{bp,raw}}^{2,\text{ev}}(i, j) \times \prod_{\ell=1,2} \Pr^{\text{ev}}[\mathcal{E}_\ell(\sigma_\ell^{\text{p}}) | \mathcal{A}_\ell] & \text{else} \end{cases} \end{aligned} \quad (18)$$

as well as  $\Pr_{\text{ss}}^{2,\text{th}}(i, s_1 \& s_2)$  in Equation (17):

$$\begin{aligned} & \mathcal{R}_{\text{ss}}^{2,\text{ev}}(i, \mathcal{A}_1 \& \mathcal{A}_2) \\ &= \begin{cases} 0 & \text{if } \exists j \text{ with } (i, j) \in \sigma_1^{\text{p}} \text{ or } (j, i) \in \sigma_1^{\text{p}} \\ 0 & \text{if } \exists j \text{ with } (i, j) \in \sigma_2^{\text{p}} \text{ or } (j, i) \in \sigma_2^{\text{p}} \\ \mathcal{R}_{\text{ss,raw}}^{2,\text{ev}}(i) \times \prod_{\ell=1,2} \Pr^{\text{ev}}[\mathcal{E}_\ell(\sigma_\ell^{\text{p}}) | \mathcal{A}_\ell] & \text{else} \end{cases} \end{aligned} \quad (19)$$

The probabilities of the partial structures  $\Pr^{\text{ev}}[\mathcal{E}_1(\sigma_1^{\text{p}}) | \mathcal{A}_1]$  and  $\Pr^{\text{ev}}[\mathcal{E}_2(\sigma_2^{\text{p}}) | \mathcal{A}_2]$  are calculated

**Input:**  
 $\mathcal{A}_1 \& \mathcal{A}_2 = \text{concatenate}(\mathcal{A}_1, \mathcal{A}_2)$   
 $C_{ss} = \dots x_1 x_1 \dots \& \dots x_2 x_2 \dots$ , where  $x$ 's are single-stranded constraints and  $x_1 \in \sigma_1^{\text{p}}, x_2 \in \sigma_2^{\text{p}}$   
**Search  $\sigma_{\text{int}}$  constrained by  $C_{ss}$ :**  
 calculate tree  $T_{\mathcal{A}_1 \& \mathcal{A}_2}$ ,  
 phylogenetic reliabilities  $\mathcal{R}_{\text{raw}}^{2,\text{ev}}$ ,  
 thermodynamic probabilities  $\Pr_{\text{raw}}^{2,\text{th}}$   
**for all  $(i, j)$  do**  
   **if  $(i, j) \in \sigma_\ell^{\text{p}}$  for  $\ell=(1,2)$  then**  
      $\mathcal{R}_{\text{bp,raw}}^{2,\text{ev}}(i, j) \leftarrow \Pr^{\text{ev}}[\mathcal{E}_\ell(\sigma_\ell^{\text{p}}) | \mathcal{A}_\ell]$   
      $\mathcal{R}_{\text{ss,raw}}^{2,\text{ev}}(i) \leftarrow 0$   
      $\Pr_{\text{bp,raw}}^{2,\text{th}}(i, j) \leftarrow \Pr^{\text{th}}[\mathcal{E}_\ell(\sigma_\ell^{\text{p}}) | s_\ell]$   
      $\Pr_{\text{ss,raw}}^{2,\text{th}}(i) \leftarrow 0$   
   **else**  
      $\mathcal{R}^{2,\text{ev}} \leftarrow \mathcal{R}_{\text{raw}}^{2,\text{ev}} \times \prod_{\ell=1,2} \Pr^{\text{ev}}[\sigma_\ell^{\text{p}} | \mathcal{A}_\ell]$   
      $\mathcal{R}^{2,\text{th}} \leftarrow \Pr_{\text{raw}}^{2,\text{th}} \times \prod_{\ell=1,2} \Pr^{\text{th}}[\sigma_\ell^{\text{p}}(s_\ell, \mathcal{A}) | s_\ell]$   
   **end if**  
    $\Rightarrow \mathcal{R}_{\text{bp}}^{2,\text{ev}}(i, j), \mathcal{R}_{\text{ss}}^{2,\text{ev}}(i) \leftarrow \text{PETfold model}$   
**end for**  
 $\sigma_{\text{int}} \leftarrow \text{MEA-structure constrained by } C_{ss}$   
**Output:**  
 $\sigma_1^{\text{p}} \cup \sigma_2^{\text{p}} \cup \sigma_{\text{int}}$

Figure 2 Pseudocode for Step 2.

as described in Equation (6). Step 2 is summarised as pseudocode in 1.

### The final scoring

To summarise the reliabilities, a combined structure will be determined using the Nussinov algorithm on the following reliabilities:

$$\mathcal{R}_{bp}^2(i, j) = \mathcal{R}_{bp}^{2, ev}(i, j, \mathcal{A}_1 \& \mathcal{A}_2) + \frac{\beta}{n} \sum_{s_1 \& s_2} \text{Pr}_{bp}^{2, th}(i, j, s_1 \& s_2)$$

$$\mathcal{R}_{ss}^2(i) = \mathcal{R}_{ss}^{2, ev}(i, \mathcal{A}_1 \& \mathcal{A}_2) + \frac{\beta}{n} \sum_{s_1 \& s_2} \text{Pr}_{ss}^{2, th}(i, s_1 \& s_2),$$

where  $\mathcal{R}_{bp}^{2, ev}(i, j, \mathcal{A}_1 \& \mathcal{A}_2)$  and  $\mathcal{R}_{ss}^{2, ev}(i, \mathcal{A}_1 \& \mathcal{A}_2)$  are defined as above,  $\text{Pr}_{bp}^{2, th}(i, j, s_1 \& s_2)$  as in Equation (16) and  $\text{Pr}_{ss}^{2, th}(i, s_1 \& s_2)$  as in Equation (17).

Note that the base pairs in  $\sigma_1^p \cup \sigma_2^p$  have a weight of 0 during folding of the constrained structure to allow for pseudoknot formation. Finally, we add the base pairs in  $\sigma_1^p \cup \sigma_2^p$  to the constrained structure of Step 2. The flow of the structure reliabilities in the pipeline is summarised in Figure 3.

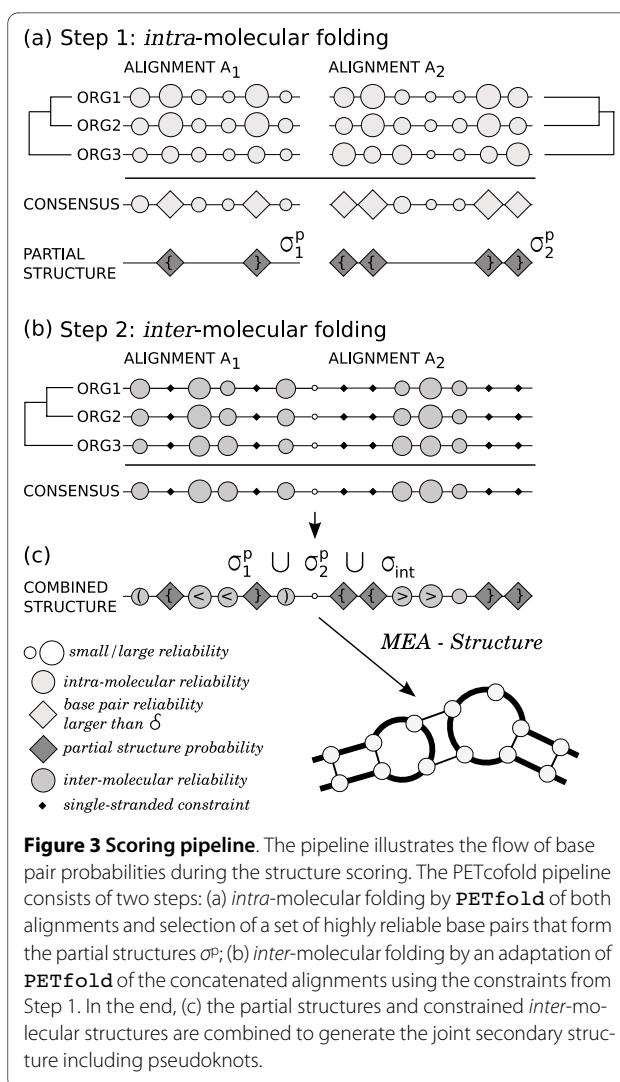
## Results and discussion

The algorithm presented herein was implemented in **PETcofold** (Seemann *et al.*, submitted). As a proof of concept, we present an example of a bacterial sRNA-mRNA interaction. The in-depth analysis is described elsewhere (Seemann *et al.*, submitted).

### Joint structure prediction of bacterial sRNA OxyS and its target mRNA fhIA

The small RNA OxyS represses the translation of the mRNA *fhIA*, which is mediated through base pairing at the ribosome binding site [11]. However, the OxyS-*fhIA* interaction involves a second binding site within the coding region of *fhIA*. Both interaction sites reside in stem loops such that OxyS and *fhIA* form a double kissing hairpin interaction.

Figure 4 shows the alignment and joint secondary structure prediction of the OxyS-*fhIA* complex, *i.e.*, the secondary structures of OxyS and *fhIA* and the interaction between them, as predicted by our algorithm. The result of the prediction without extending the constrained stems is shown in Figure 4a, and the result with the extension of the constrained stems is shown in Figure 4b.



**Figure 3 Scoring pipeline.** The pipeline illustrates the flow of base pair probabilities during the structure scoring. The PETcofold pipeline consists of two steps: (a) *intra*-molecular folding by **PETfold** of both alignments and selection of a set of highly reliable base pairs that form the partial structures  $\sigma^p$ ; (b) *inter*-molecular folding by an adaptation of **PETfold** of the concatenated alignments using the constraints from Step 1. In the end, (c) the partial structures and constrained *inter*-molecular structures are combined to generate the joint secondary structure including pseudoknots.

For OxyS-*fhIA*, our algorithm was able to consistently predict one of the two interaction sites. The second interaction site, which is situated in the *fhIA* coding region, was only predicted when the constrained stems were not extended in Step 1 of our algorithm. Otherwise the stem of *fhIA* that resides the second interaction site was extended both by inner and outer base pairs. Consequently, the unpaired region of the hairpin containing the second interaction site became shorter such that no interaction was predicted at this site.

### Algorithmic restrictions and potentials

The algorithm supports pseudoknots between the *intra*-molecular and *inter*-molecular base pairs, while the time complexity of  $O(N \times I \times L^3)$  is much lower than that of other approaches with the same ability. The time complexity is in the magnitude of **PET-fold** for the added sequence length  $L$  of both alignments, and it is linear with respect to the number of sequences  $N$  in the alignments



to an unreasonable level. Furthermore, the energy contribution of the cofolding step might be slightly biased due to the constraint of the partial structures as single-stranded. We partly solve the resulting *intra*-molecular false predictions by extending the reliable stems in the partial structures, and, as already mentioned above, the RNACoFold algorithm and scoring scheme could be adapted to handle base pair constraints as single-stranded.

Furthermore, there is a limitation of the presented method with regard to interaction sites that are located outside conserved RNA structures. These regions are hard to align if they are, in addition, sequentially unconserved. Thus, our method will most likely miss binding sites located in unstructured and otherwise unrelated regions, *e.g.*, miRNA target sites in UTR regions. However, once a correct alignment is found for these regions, then the presented approach still works if the interaction region is conserved or shows enough covariation.

Our algorithm is able to predict pseudoknots between the *intra*-molecular and (inter-)molecular base pairs. In addition, we are interested in more pseudoknots that can be predicted in a similar way using a pipeline of constrained structures. During an iteration of Step 2, additional reliable partial *inter*-molecular structures are constrained as long as new reliable base pairs appear. The final consensus structure is the union of all cofolding base pairs and the partial structures. The main unsolved problem is the weighted combination of the decreasing partial structure probabilities in one scoring scheme when the amount of constraints increases with each iteration.

## Conclusions

In summary, we introduced an extension of the PET-fold algorithm for the identification of interactions between two sets of multiple aligned RNA sequences, which exploits compensating base changes while taking *intra*-molecular partial structures and interaction sites into account. The implementation of the algorithm in PETcofold and its application are described in Seemann *et al.* (submitted).

## Additional material

**Additional file 1 Probability distributions of the Pfold model and Implications of the independence property.** Combined probability distributions of the Pfold model and the implications of the independence property (Equation (11)) for partial structures are described.

### Competing interests

The authors declare that they have no competing interests.

### Authors' contributions

SES, RB and JG designed the algorithm. SES implemented the algorithm. ASR designed and performed the analysis of the algorithm. RB drafted the manuscript. All authors contributed to the manuscript, read and approved the final manuscript.

## Acknowledgements

We thank Bjarne Knudsen for inspiring discussions about extensions of the Pfold method.

This work was supported by the Lundbeck Foundation (grant 374/06 to J.G.), the Danish Research Council for technology and production (grant 274-09-0282 to J.G.), the Danish Center for Scientific Computation (J.G.), the German Federal Ministry of Education and Research (BMBF grant 0313921 FRISYS to R.B.) and the German Research Foundation (DFG grant BA 2168/2-1 SPP 1258 to A.S.R. and R.B.).

## Author Details

<sup>1</sup>Center for non-coding RNA in Technology and Health, IBHV, University of Copenhagen, Grønnegårdsvej 3, Frederiksberg C, 1870, Denmark and

<sup>2</sup>Bioinformatics Group, University of Freiburg, Georges-Köhler-Allee 106, Freiburg, 79110, Germany

Received: 10 April 2010 Accepted: 21 May 2010

Published: 21 May 2010

## References

1. Washietl S, Hofacker IL, Lukasser M, Hüttenhofer A, Stadler PF: **Genome-wide mapping of conserved RNA secondary structure structures predicts thousands of functional non-coding RNAs in human.** *Nature Biotechnology* 2005, **23**:1383-1390.
2. Pedersen JS, Bejerano G, Siepel A, Rosenbloom K, Lindblad-Toh K, Lander ES, Kent J, Miller W, Haussler D: **Identification and classification of conserved RNA secondary structures in the human genome.** *PLoS Comput Biol* 2006, **2**:e33.
3. Torarinsson E, Sawera M, Havgaard JH, Fredholm M, Gorodkin J: **Thousands of corresponding human and mouse genomic regions unalignable in primary sequence contain common RNA structure.** *Genome Research* 2006, **16**:885-889. [Erratum in: *Genome Res.* 2006 16:1439].
4. Uzilov AV, Keegan JM, Mathews DH: **Detection of non-coding RNAs on the basis of predicted secondary structure formation free energy change.** *BMC Bioinformatics* 2006, **7**:173.
5. Washietl S, Pedersen JS, Korb J, Stocsits C, Gruber AR, Hackermüller J, Hertel J, Lindemeyer M, Reiche K, Tanzer A, Ucla C, Wyss C, Antonarakis SE, Denoeud F, Lagarde J, Drenkow J, Kapranov P, Gingeras TR, Guigó R, Snyder M, Gerstein MB, Reymond A, Hofacker IL, Stadler PF: **Structured RNAs in the ENCODE selected regions of the human genome.** *Genome Research* 2007, **17**:852-864.
6. Will S, Reiche K, Hofacker IL, Stadler PF, Backofen R: **Inferring noncoding RNA families and classes by means of genome-scale structure-based clustering.** *PLoS Computational Biology* 2007, **3**:e65.
7. Torarinsson E, Yao Z, Wiklund ED, Bramsen JB, Hansen C, Kjems J, Tommerup N, Ruzzo WL, Gorodkin J: **Comparative genomics beyond sequence based alignments: RNA structures in the ENCODE regions.** *Genome Research* 2008, **18**:242-251.
8. Backofen R, Hess WR: **Computational prediction of sRNAs and their targets in bacteria.** *RNA Biol* 2010, **7**.
9. Mückstein U, Tafer H, Hackermüller J, Bernhart SH, Stadler PF, Hofacker IL: **Thermodynamics of RNA-RNA binding.** *Bioinformatics* 2006, **22**(10):1177-82.
10. Busch A, Richter AS, Backofen R: **IntaRNA: efficient prediction of bacterial sRNA targets incorporating target site accessibility and seed regions.** *Bioinformatics* 2008, **24**(24):2849-56.
11. Argaman L, Altuvia S: **fhfA repression by OxyS RNA: kissing complex formation at two sites results in a stable antisense-target RNA complex.** *Journal of Molecular Biology* 2000, **300**(5):1101-12.
12. Andronescu M, Zhang ZC, Condon A: **Secondary structure prediction of interacting RNA molecules.** *Journal of Molecular Biology* 2005, **345**(5):987-1001.
13. Bernhart SH, Tafer H, Mückstein U, Flamm C, Stadler PF, Hofacker IL: **Partition function and base pairing probabilities of RNA heterodimers.** *Algorithms Mol Biol* 2006, **1**:3.
14. Dirks RM, Bois JS, Schaeffer JM, Winfree E, Pierce NA: **Thermodynamic Analysis of Interacting Nucleic Acid Strands.** *SIAM Review* 2007, **49**:65-88.
15. Zuker M: **Prediction of RNA secondary structure by energy minimization.** *Methods in Molecular Biology* 1994, **25**:267-94.



16. Hofacker IL, Fontana W, Stadler PF, Bonhoeffer S, Tacker M, Schuster P: **Fast Folding and Comparison of RNA Secondary Structures.** *Monatshette Chemie* 1994, **125**:167-188.
17. Pervouchine DD: **IRIS: intermolecular RNA interaction search.** *Genome Inform* 2004, **15**(2):92-101.
18. Alkan C, Karakoc E, Nadeau JH, Sahinalp SC, Zhang K: **RNA-RNA interaction prediction and antisense RNA target search.** *Journal of Computational Biology* 2006, **13**(2):267-82.
19. Chitsaz H, Salari R, Sahinalp SC, Backofen R: **A partition function algorithm for interacting nucleic acid strands.** *Bioinformatics* 2009, **25**(12):1365-73.
20. Huang FWD, Qin J, Reidys CM, Stadler PF: **Partition function and base pairing probabilities for RNA-RNA interaction prediction.** *Bioinformatics* 2009, **25**(20):2646-54.
21. Huang FWD, Qin J, Reidys CM, Stadler PF: **Target prediction and a statistical sampling algorithm for RNA-RNA interaction.** *Bioinformatics* 2010, **26**(2):175-81.
22. Chitsaz H, Backofen R, Sahinalp SC: **biRNA: Fast RNA-RNA Binding Sites Prediction.** *Proc. of the 9th Workshop on Algorithms in Bioinformatics (WABI), Volume 5724 of Lecture Notes in Computer Science* 2009:25-36.
23. Salari R, Backofen R, Sahinalp SC: **Fast prediction of RNA-RNA Interaction.** *Proc. of the 9th Workshop on Algorithms in Bioinformatics (WABI), Volume 5724 of Lecture Notes in Computer Science* 2009:261-272.
24. Salari R, Möhl M, Will S, Sahinalp SC, Backofen R: **Time and space efficient RNA-RNA interaction prediction via sparse folding.** *Proc of RECOMB 2010* 2010 in press.
25. Seemann SE, Gorodkin J, Backofen R: **Unifying evolutionary and thermodynamic information for RNA folding of multiple alignments.** *Nucleic Acids Research* 2008, **36**:6355-6362.
26. Gaspin C, Westhof E: **An interactive framework for RNA secondary structure prediction with a dynamical treatment of constraints.** *J Mol Biol* 1995, **254**:163-174.
27. Jabbari H, Condon A, Pop A, Pop C, Zhao Y: **HFold: RNA Pseudoknotted Secondary Structure Prediction Using Hierarchical Folding.** In *Algorithms in Bioinformatics, 7th International Workshop, WABI Philadelphia, PA, USA, September 8-9, 2007, Proceedings* 2007:323-334.
28. Knudsen B, Hein JJ: **RNA secondary structure prediction using stochastic context-free grammars and evolutionary history.** *Bioinformatics* 1999, **15**:446-454.
29. McCaskill JS: **The equilibrium partition function and base pair binding probabilities for RNA secondary structure.** *Biopolymers* 1990, **29**(6-7):1105-19.
30. Nussinov R, Pieczenik G, Griggs JR, Kleitman DJ: **Algorithms for Loop Matchings.** *SIAM Journal on Applied Mathematics* 1978, **35**:68-82.
31. Ding Y, Chan CY, Lawrence CE: **RNA secondary structure prediction by centroids in a Boltzmann weighted ensemble.** *RNA* 2005, **11**(8):1157-66.
32. Knudsen B, Andersen ES, Damgaard C, Kjems J, Gorodkin J: **Evolutionary rate variation and RNA secondary structure prediction.** *Comput Biol Chem* 2004, **28**(3):219-226.
33. Waterhouse AM, Procter JB, Martin DMA, Clamp M, Barton GJ: **Jalview Version 2-a multiple sequence alignment editor and analysis workbench.** *Bioinformatics* 2009, **25**(9):1189-91.

doi: 10.1186/1748-7188-5-22

**Cite this article as:** Seemann et al., Hierarchical folding of multiple sequence alignments for the prediction of structures and RNA-RNA interactions *Algorithms for Molecular Biology* 2010, **5**:22

**Submit your next manuscript to BioMed Central and take full advantage of:**

- Convenient online submission
- Thorough peer review
- No space constraints or color figure charges
- Immediate publication on acceptance
- Inclusion in PubMed, CAS, Scopus and Google Scholar
- Research which is freely available for redistribution

Submit your manuscript at  
www.biomedcentral.com/submit

