

## Topology and prediction of RNA pseudoknots

Reidys, Christian; Huang, Fenix W.D.; Ellegaard Andersen, Jørgen; Penner, Robert; Stadler, Peter F.; Nebel, Markus E.

*Published in:*  
Bioinformatics

*DOI:*  
10.1093/bioinformatics/btr090

*Publication date:*  
2011

*Document version:*  
Submitted manuscript

*Citation for published version (APA):*  
Reidys, C., Huang, F. W. D., Ellegaard Andersen, J., Penner, R., Stadler, P. F., & Nebel, M. E. (2011). Topology and prediction of RNA pseudoknots. *Bioinformatics*, 27(8). <https://doi.org/10.1093/bioinformatics/btr090>

Go to publication entry in University of Southern Denmark's Research Portal

### Terms of use

This work is brought to you by the University of Southern Denmark.  
Unless otherwise specified it has been shared according to the terms for self-archiving.  
If no other license is stated, these terms apply:

- You may download this work for personal use only.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying this open access version

If you believe that this document breaches copyright please contact us providing details and we will investigate your claim.  
Please direct all enquiries to [puresupport@bib.sdu.dk](mailto:puresupport@bib.sdu.dk)

# Topology and prediction of RNA pseudoknots

Christian M. Reidys<sup>1,2\*</sup>, Fenix W.D. Huang<sup>1</sup>, Jørgen E. Andersen<sup>3</sup>, Robert C. Penner<sup>3,4</sup>, Peter F. Stadler<sup>5–10</sup>, and Markus E. Nebel<sup>11</sup>

<sup>1</sup>Center for Combinatorics, LPMC-TJKLC, Nankai University Tianjin 300071, P.R. China

<sup>2</sup>College of Life Science, Nankai University Tianjin 300071, P.R. China

<sup>3</sup>Center for Quantum Geometry of Moduli Spaces Aarhus University, DK-8000 Århus C, Denmark

<sup>4</sup>Math and Physics Departments, California Institute of Technology, Pasadena, California, USA

<sup>5</sup>Bioinformatics Group, Department of Computer Science, and Interdisciplinary Center for Bioinformatics, University of Leipzig, Härtelstraße 16-18, D-04107 Leipzig, Germany.

<sup>6</sup>Max Planck Institute for Mathematics in the Sciences, Inselstraße 22, D-04103 Leipzig, Germany

<sup>7</sup>RNomics Group, Fraunhofer IZI, Perlickstraße 1, D-04103 Leipzig, Germany

<sup>8</sup>Inst. f. Theoretical Chemistry, University of Vienna, Währingerstraße 17, A-1090 Vienna, Austria

<sup>9</sup>Center for non-coding RNA in Technology and Health, University of Copenhagen, Grønnegårdsvej 3, DK-1870 Frederiksberg, Denmark

<sup>10</sup>The Santa Fe Institute, 1399 Hyde Park Rd., Santa Fe, New Mexico, USA

<sup>11</sup>Department of Computer Science, University of Kaiserslautern, Germany

Received on \*\*\*\*; revised on \*\*\*\*; accepted on \*\*\*\*

Associate Editor: \*\*\*\*

## ABSTRACT

**Motivation:** Several dynamic programming algorithms for predicting RNA structures with pseudoknots have been proposed that differ dramatically from one another in the classes of structures considered.

**Results:** Here we use the natural topological classification of RNA structures in terms of irreducible components that are embeddable in surfaces of fixed genus. We add to the conventional secondary structures four building blocks of genus one in order to construct certain structures of arbitrarily high genus. A corresponding unambiguous multiple context free grammar provides an efficient dynamic programming approach for energy minimization, partition function, and stochastic sampling. It admits a topology-dependent parameterization of pseudoknot penalties that increases the sensitivity and positive predictive value of predicted base pairs by 10–20% compared to earlier approaches. More general models based on building blocks of higher genus are also discussed.

**Availability:** The source code of `gfold` is freely available at <http://www.combinatorics.cn/cbpc/gfold.tar.gz>

**Contact:** [duck@santafe.edu](mailto:duck@santafe.edu)

**Supplementary information:** Supplementary material containing a complete presentation of the algorithms, full proofs of theorems, and detailed performance data are available at *Bioinformatics online*.

## 1 INTRODUCTION

The global conformation of RNA molecules is to a large extent determined by topological constraints encoded at the level of secondary structure, i.e., by the mutual arrangements of the base paired

helices (Bailor *et al.*, 2010). In this context, secondary structure is understood in a wider sense that includes pseudoknots. Although the vast majority of RNAs has simple, i.e., pseudoknot-free, secondary structure, `PseudoBase` (Taufers *et al.*, 2009) lists more than 250 records of pseudoknots determined by a variety of experimental and computational techniques including crystallography, NMR, mutational experiments, and comparative sequence analysis. In many cases, they are crucial for molecular function. Examples include the catalytic cores of several ribozymes (Doudna and Cech, 2002), programmed frameshifting (Namy *et al.*, 2006), and telomerase activity (Theimer *et al.*, 2005), reviewed in (Staple and Butcher, 2005; Giedroc and Cornish, 2009).

Secondary structures can be interpreted as matchings in a graph of permissible base pairs (Tabaska *et al.*, 1998). The energy of RNA folding is dominated by the stacking of adjacent base pairs, not by the hydrogen bonds of the individual base pairs (Mathews *et al.*, 1999). In contrast to maximum weighted matching, the general RNA folding problem with a stacking-based energy function is NP-complete (Akutsu, 2000; Lyngsø and Pedersen, 2000). The most commonly used RNA secondary structure prediction tools, including `mfold` (Zuker, 1989) and the `Vienna RNA Package` (Hofacker *et al.*, 1994), therefore exclude pseudoknots.

Polynomial-time dynamic programming (DP) algorithms can be devised, however, for certain restricted classes of pseudoknots. In contrast to the  $O(N^2)$  space and  $O(N^3)$  time solution for simple secondary structures (Waterman, 1978; Nussinov *et al.*, 1978; Zuker and Stiegler, 1981), however, most of these approaches are computationally much more demanding. The design of pseudoknot folding algorithms thus has been governed more by the need to limit computational cost and achieve a manageable complexity of the recursion than the conscious choice of a particularly natural search space of RNA structures. As a case in point, the class

\*to whom correspondence should be addressed. Phone: \*86-22-2350-6800; Fax: \*86-22-2350-9272; [duck@santafe.edu](mailto:duck@santafe.edu)

of structures underlying the algorithm by Rivas and Eddy (1999) (R&E-structures, `pknot-R&E`) was characterized only in a subsequent publication (Rivas and Eddy, 2000). The following references provide a certainly incomplete list of DP approaches to RNA structure prediction using different structure classes characterized in terms of recursion equations and/or stochastic grammars: Rivas and Eddy (1999); Uemura Y. *et al.* (1999); Akutsu (2000); Lyngsø and Pedersen (2000); Cai *et al.* (2003); Dirks and Pierce (2003); Deogun *et al.* (2004); Reeder and Giegerich (2004); Li and Zhu (2005); Matsui *et al.* (2005); Kato *et al.* (2006); Chen *et al.* (2009). The inter-relationships of some of these classes of RNA structures have been clarified in part by Condon *et al.* (2004) and Rødland (2006). In addition to these exact algorithms, a plethora of heuristic approaches to pseudoknot prediction have been proposed in the literature; see e.g., (Metzler and Nebel, 2008; Chen, 2008) and the references therein.

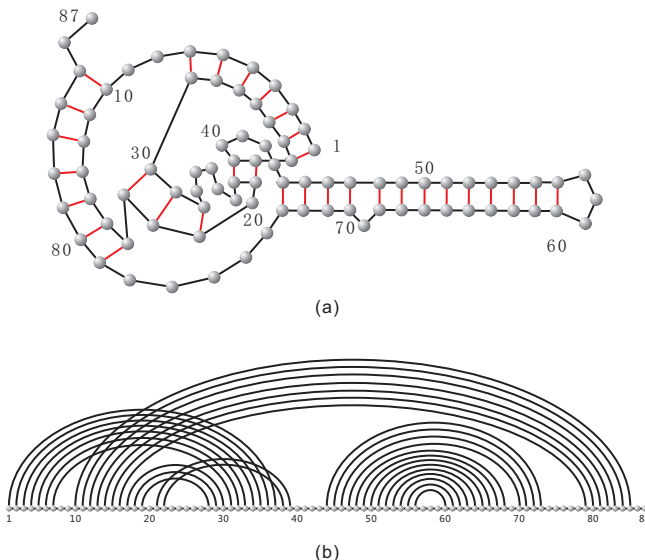
At least three distinct classification schemes of RNA contact structures have been proposed: Haslinger and Stadler (1999) suggested using book-embeddings, Jin *et al.* (2008) focused on the maximal set of pairwise crossing base pairs, and Bon *et al.* (2008) based the classification on topological embeddings. While these classifications have in common that simple secondary structure forms the most primitive class of structures, they differ already in the construction of the first non-trivial class of pseudoknots. Despite their mathematical appeal, however, no efficient (polynomial-time) algorithms are available for predicting pseudoknotted structures even in the simplest case of 3-noncrossing RNA structures. A practically workable approach to 3-noncrossing structures requires the enumeration of an exponentially growing number of diagrams which are then “filled in” by means of DP (Huang *et al.*, 2009); a Monte-Carlo approach utilizing the topological approach with a very simple matching-like energy model was explored by (Vernizzi and Orland, 2005).

In this contribution, we show that the topological classification of RNA structures can be translated into efficient DP algorithms. To this end, we introduce  $\gamma$ -structures and prove that they can be derived from a *finite* family of abstract shapes called shadows. In Theorem 2.3, we enumerate these four shadows for  $\gamma = 1$ , which can be cast as explicit construction rules for a unique multiple context-free grammar (Section 2.3). Corresponding DP algorithms for energy minimization, partition function, and Boltzmann-sampling functionalities are implemented in the software package `gfold`. An important feature is that  $\gamma$ -structures can be treated algorithmically like pseudoknot-free secondary structures in sense that there are finitely many motifs, i.e., shadows, for fixed  $\gamma$ , each of which is assigned a specific energy. Because of the multiplicity of motifs, which rapidly increases with  $\gamma$ , this allows for a more detailed energy model of pseudoknotted structures based on their topological complexity.

## 2 RESULTS

### 2.1 Topology of RNA Structures

*Diagram Representation.* RNA molecules are linear biopolymers consisting of the four nucleotides **A**, **U**, **C**, and **G** characterized by a sequence endowed with a unique orientation (5' to 3'). Each nucleotide can interact (base pair) with at most one other nucleotide by means of specific hydrogen bonds. Only the Watson-Crick pairs **GC**



**Fig. 1.** RNA structure as planar graph (hydrogen bonds (resp. backbone) represented by red (resp. black) edges) and diagram.

and **AU** as well as the wobble **GU** are admissible. These base pairs determine the secondary structure. Note that we have neglected here base triples and other types of more complex interactions. Secondary structures can thus be represented as graphs where nucleotides are represented by vertices, the backbone of the molecule as well as the hydrogen bonds are represented by edges; see Fig. 1 (a). More conveniently, we use the convention to represent the backbone of the polymer by a horizontally drawn chain. As before, this chain consists of vertices and arcs respectively representing the nucleotides and covalent bonds. However, the edges representing the base pairs now are depicted as arcs in the upper half-plane; see Fig. 1 (b). We call this representation the *diagram* of the molecule.

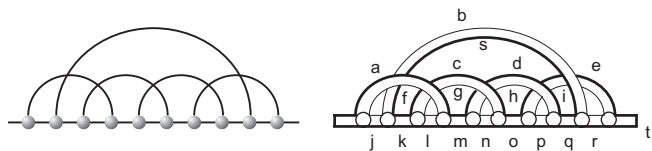
Thus, we shall identify a structure with a labelled graph over the vertex set  $[N] = \{1, 2, \dots, N\}$  represented by drawing the vertices  $1, 2, \dots, N$  on a horizontal line in the natural order and the arcs  $(i, j)$ , where  $i < j$ , in the upper half-plane.

*Fatgraph representation.* In order to understand the topological properties of RNA molecules we need to pass from the picture of RNA as diagrams or contact-graphs to that of topological surfaces. Only the associated surface carries the important invariants leading to a meaningful filtration of RNA structures. Formally, we will view an RNA molecule as a topological surface (Andersen *et al.*, 2010). The main idea is to “thicken” the edges into (untwisted) bands or ribbons and to expand each vertex to a disk as shown in Fig. 2. This inflation of edges leads to a fatgraph  $\mathbb{D}$  (Loebl and Moffatt, 2008; Penner *et al.*, 2010).

A fatgraph, sometimes also called “ribbon graph” or “map”, is a graph equipped with a cyclic ordering of the incident half-edges at each vertex. Thus,  $\mathbb{D}$  refines its underlying graph  $D$  insofar as it encodes the ordering of the ribbons incident on its disks. In the



**Fig. 2.** Inflation of edges and vertices to ribbons and disks. Here we have four vertices, five edges and one boundary component ( $\vec{a}, \vec{b}, \vec{c}, \vec{d}, \vec{e}, \vec{f}, \vec{g}, \vec{h}, \vec{i}, \vec{j}$ ). The corresponding surface has Euler characteristic  $\chi = v - e + r = 0$  and genus  $g = 1$ , see eqs (2.1) and (2.2).



**Fig. 3.** Computing the number of boundary components. The diagram contains  $5 + 9$  edges and 10 vertices. We follow the alternating paths described in the text and observe that there are exactly two boundary components (bold and thin). According to eq. (2.1), the genus of the diagram is given by  $1 - \frac{1}{2}(10 - 14 + 2) = 2$ , see SM, Fig. S6 for details.

following we will deal with orientable ribbon graphs<sup>1</sup>. Each ribbon has two boundaries. The first one in counterclockwise order is labeled by an arrowhead, see Fig. 2. A  $\mathbb{D}$ -cycle or  $\mathbb{D}$ -boundary component is then constructed by following these directed boundaries from disk to disk thereby alternating between base pair ribbons and backbone, with the exception of the segment of the boundary component that travels along the bottom of the backbone using only backbone bonds, as shown in Figs. 2 and 3. We give a brief tutorial on how to compute boundary components in the SM, Fig. S6. Topological invariants such as the number of boundary components of the fatgraph  $\mathbb{D}$  can thus be computed directly from the underlying diagram  $D$ . Furthermore, fatgraphs can be succinctly stored and conveniently manipulated on the computer as pairs of permutations (Penner *et al.*, 2010).

The fatgraph  $\mathbb{D}$  gives rise to a unique surface  $X_{\mathbb{D}}$ , and each  $\mathbb{D}$ -cycle corresponds to a boundary component of  $X_{\mathbb{D}}$ , whose Euler characteristic and genus are given by

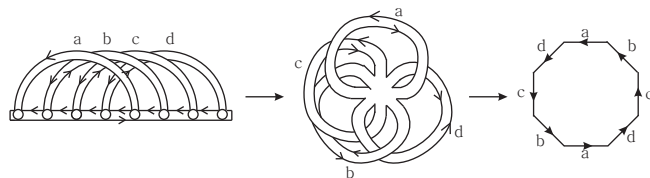
$$\chi(X_{\mathbb{D}}) = v - e + r \quad (2.1)$$

$$g(X_{\mathbb{D}}) = 1 - \frac{1}{2}\chi(X_{\mathbb{D}}), \quad (2.2)$$

where  $v, e, r$  denotes the number of discs, ribbons and boundary components in  $\mathbb{D}$  (Massey, 1967). The graph  $D$  can readily be obtained by continuously contracting the ribbons and discs of  $\mathbb{D}$ .

We next make use of an additional feature of RNA structures, namely, that the backbone forms a unique oriented chain determined by the covalent bonds. Thus, the backbone can be collapsed to a single disk since the surface is orientable: in absence of twisted ribbons, there is no particular information in the backbone itself. Indeed, the procedure can be undone by re-inflating the disk and rebuilding the backbone. The contraction of the  $N$  vertices to a single one and

<sup>1</sup> ribbons may also be allowed to twist giving rise to possibly non-orientable surfaces (Massey, 1967)



**Fig. 4.** Reduction to fatgraphs with a single vertex. Contracting the backbone of a diagram into a single vertex decreases the length of the boundary components and preserves the genus. The contracted fatgraph is equivalent to the labeled directed cycle. The backbone of the polymer can be recovered by re-inflating the disk into the backbone. The polygon (r.h.s.) represents the standard 2D-model of a surface as discussed in (Massey, 1967).

the removal of the  $(N - 1)$  covalent bonds therefore preserves the Euler characteristic and genus, see Fig. 4.

Using the collapsed fatgraph<sup>2</sup> we see that the relation between the genus of the surface and the number of boundary components is determined by the number of arcs in the upper half-plane, namely,

$$2 - 2g - r = 1 - n, \quad (2.3)$$

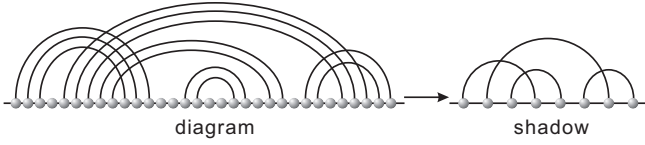
where  $n$  is number of base pairs and  $r$  the number of boundary components. The latter can be computed easily and therefore controls the genus of the molecules. Eq. (2.3) follows from eqns. (2.2) and (2.1), which together yield  $2 - 2g - r = v - e$ , and the observation that the contracted graph has  $e = n$  arcs and a single ( $v = 1$ ) vertex.

## 2.2 $\gamma$ -structures

The *shadow* of a diagram (RNA structure) is obtained by removing all noncrossing arcs, collapsing all isolated vertices and replacing all remaining stacks (i.e., adjacent parallel arcs) by single arcs; see Fig. 5. Shadows can be seen as a generalization of shape abstractions (Giegerich *et al.*, 2004) to pseudoknotted structures (Reidys and Wang, 2010). Similar to the process of contracting the backbone into a single vertex, the projection into a shadow changes neither genus nor the number of boundary components (Andersen *et al.*, 2010). All information on stack-lengths and on noncrossing components of the structure is lost in the process however. We shall see that the set of structures with shadow  $\mathfrak{S}$  can nevertheless be reconstructed efficiently. To this end we will show that, for fixed genus  $g$ , there are only *finitely many* distinct shadows  $S_g$ , which will play a central role in constructing folding algorithms.

A diagram is *irreducible* (or connected) (Kleitman, 1970) if for any two arcs there is a sequence of arcs so that consecutive arcs cross one other. A shadow is not necessarily irreducible but may be composed of multiple irreducible components or blocks, see Fig. 6 (1). Any shadow (and in general, any diagram) can be decomposed iteratively by removing irreducible components from bottom to top, i.e., so that that there is no component “inside” the one just removed. Note that the set  $\mathbf{I}_{\mathfrak{S}}$  of irreducible components of the set of shadows,  $\mathfrak{S}(S)$ , equals the set of shadows of the irreducible components of

<sup>2</sup> in order to relate this to the standard 2D-models of surfaces derived from triangulations: from the collapsed fatgraph we can derive the *polygonal model of the surface*  $X_{\mathbb{D}}$ , i.e., a  $2n$ -gon in which edges are identified in pairs; see Fig. 4



**Fig. 5.** The shadow of a diagram is obtained by removing all noncrossing arcs and isolated vertices and collapsing all resulting stacks into single arcs. While taking shadows is a significant reduction, the key topological invariants of genus and number of boundary components remain invariant.



**Fig. 6.**  $\gamma$ -structures: we display the shadow of a 1-structure (left) having topological genus two and the shadow of the HDV-structure (right) (Ferré-D’Amaré *et al.*, 1998), a 2-structure having also genus two. Although both shadows have genus two, the HDV structure cannot be generated iteratively via successive removals of  $S_1$ -elements and stacked arcs. The structure displayed on the left is derived via two  $S_1$ -substructures.

the diagram  $S$ . Furthermore, the genus of  $\mathfrak{S}(S)$  is the sum of the genera of its irreducible components, i.e.,

$$g(S) = g(\mathfrak{S}(S)) = \sum_{\mathfrak{S}' \in \mathbf{I}_{\mathfrak{S}(S)}} g(\mathfrak{S}'). \quad (2.4)$$

It seems natural, therefore, to determine the complexity of a structure by the maximal genus of the components of its shadows. More precisely, we say that  $S$  is a  $\gamma$ -structure if  $g(\mathfrak{S}') \leq \gamma$  holds for all irreducible components of the shadows  $\mathfrak{S}(S)$ . By definition, a  $\gamma$ -structure can thus be constructed from the set  $S_\gamma$  of shadows of genus at most  $\gamma$  by inserting certain noncrossing arcs, see Fig. 6. The simplest class of structures are of course 0-structures, obtained by placing noncrossing arcs over the empty structure.

**LEMMA 2.1.** *An RNA structure is a 0-structure if and only if it is a simple secondary structure. In particular, a 0-structure always has genus  $g = 0$ .*

**PROOF.** We first observe that a diagram of genus zero contains no crossing arcs. This follows from the fact that genus is a monotone non-decreasing function of the number of arcs (see eq. (2.3)) and that the genus of the matching (H) consisting of two mutually crossing arcs has only one boundary component and hence genus one; see Fig. 2. Second, we observe by induction on the number of arcs that each new noncrossing arc contributes a new boundary component and  $2 - 2g - (r + 1) = 1 - (n + 1)$  shows that the genus remains zero. Structures consisting only of noncrossing arcs therefore have genus zero.

Next, we consider structures of arbitrary genus. For their analysis, diagrams without isolated points, i.e., matchings, play a central role. Let  $\mathcal{C}_g(n)$  be the set of matchings of genus  $g$  with  $n$  arcs, and let  $\mathbf{c}_g(n) := |\mathcal{C}_g(n)|$  denote its cardinality. As shown by Andersen *et al.* (2010), the generating function  $\mathbf{C}_g(z) = \sum_{n \geq 0} \mathbf{c}_g(n) z^n$  is

given by

$$\mathbf{C}_g(z) = P_g(z) \frac{\sqrt{1-4z}}{(1-4z)^{3g}}, \quad g \geq 1, \quad (2.5)$$

where  $P_g(z)$  is an integral polynomial of degree  $(3g - 1)$  such that  $P_g(1/4) \neq 0$ . The number of genus zero matchings are well-known to be given by the Catalan numbers, and eq. (2.5) allows the derivation of explicit formulas for higher genera, for instance,

$$\mathbf{c}_1(n) = \frac{2^{n-2}(2n-1)!!}{3(n-2)!}, \quad \mathbf{c}_2(n) = \frac{2^{n-4}(5n-2)(2n-1)!!}{90(n-4)!}.$$

Furthermore, the number  $\mathbf{c}_g(2g)$  of matchings of genus  $g$  having exactly  $2g$  arcs, i.e., matchings having exactly one boundary component, is the coefficient of  $z^{2g}$  in  $P_g(z)$  and is given by

$$\mathbf{c}_g(2g) = \frac{(4g)!}{4^g(2g+1)!}. \quad (2.6)$$

Explicitly, we have  $\mathbf{c}_1(2) = 1$ ,  $\mathbf{c}_2(4) = 21$  and  $\mathbf{c}_3(6) = 1485$  for example. These particular matchings will serve as “seeds” for our folding algorithm. More precisely, we shall use the following:

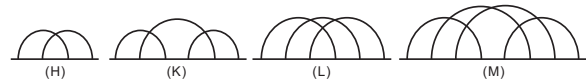
**THEOREM 2.2.** *For arbitrary genus  $g$ , the set  $S_g$  of shadows is finite. Every shadow in  $S_g$  contains at least  $2g$  and at most  $(6g - 2)$  arcs.*

The special case  $g = 1$ , on which we focus in the algorithmic part of this contribution, is explicated in the Supplementary Material (SM).

**PROOF.** First note that if there is more than one boundary component, then there must be an arc with different boundary components on its two sides, and removing this arc decreases  $r$  by exactly one while preserving  $g$  since the number of arcs is given by  $n = 2g + r - 1$ . Furthermore, if there are  $\nu_\ell$  boundary components of length  $\ell$  in the polygonal model, then  $2n = \sum_\ell \ell \nu_\ell$  since each side of each arc is traversed once by the boundary. For a shadow,  $\nu_1 = 0$  by definition, and  $\nu_2 \leq 1$  as one sees directly. It therefore follows that  $2n = \sum_\ell \ell \nu_\ell \geq 3(r-1) + 2$ , so  $2n = 4g + 2r - 2 \geq 3r - 1$ , i.e.,  $4g - 1 \geq r$ . Thus we have  $n = 2g + (4g - 1) - 1 = 6g - 2$ , i.e. any shadow can contain at most  $6g - 2$  arcs. The lower bound  $2g$  follows directly from  $n = 2g + r - 1$  by observing  $r = 1$ .

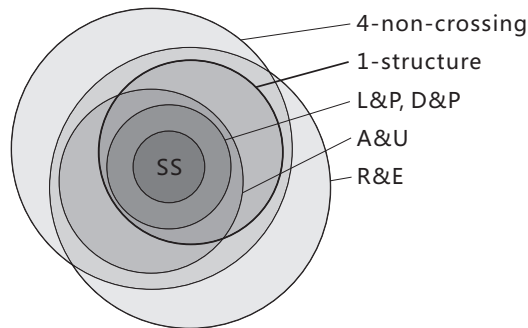
Many  $S_g$ -shadows are in fact  $\gamma$  structures for some  $\gamma < g$ , that is, they can be constructed from elements of  $S_\gamma$ . One key result of this contribution is the following characterization of 1-structures:

**THEOREM 2.3.** *An RNA structure is a 1-structure if and only if its shadow can be decomposed by iteratively removing one of the four shadows*



*In particular, 1-structures can have arbitrarily large topological genus.*

**PROOF.** We only give a sketch here and refer to the SM for a full proof. First, we observe that taking the shadow preserves genus. Since (H) is the unique matching with two arcs of genus  $g = 1$ , it is contained in every matching of genus  $g = 1$ . An arc crossing



**Fig. 7.** Venn diagram of important classes of structures with pseudoknots. The mutual relationships of pseudoknot-free secondary structure (SS), the two H-shadow classes D&P and L&P, and the classes A&U and R&E, resp., were already described by Condon *et al.* (2004). 1-structures and 4-noncrossing structures are added here.

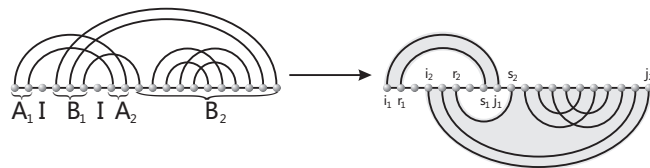
into (H) preserves the genus and leads to either (K) or (L). While every arc added to (K) increases the genus, there is one possibility to preserve the genus when adding an arc to (L), namely, the addition leading to (M). It remains to observe that no further arc can be added to (M).

Before proceeding to algorithmic considerations we briefly compare the class of  $\gamma$ -structures with other classes of pseudoknots. Condon *et al.* (2004) investigated the structure classes L&P (Lyngsø and Pedersen, 2000), D&P (Dirks and Pierce, 2003), A&U (Akutsu, 2000), and R&E-class (Rivas and Eddy, 1999). The L&P- and D&P-class are based on the H-type shadow depicted in Theorem 2.3 and hence are proper subsets of the 1-structures. The A&U-class does not cover shadow M but on the other hand contains some configurations that are not 1-structures, and even the 2-structures do not completely contain the A&U-class. Nevertheless, the A&U-class is small: there are more 1-structures than A&U-structures for any given sequence length (Nebel and Weinberg, 2011).

The R&E class does not impose a limit on the genus of the shadow and hence contains  $\gamma$ -structure with arbitrarily large  $\gamma$ . Conversely, Fig. 3 shows a 2-structure that is not contained in the R&E class. This example is minimal, i.e., all 1-structures are contained in R&E. Similarly, the set of  $k$ -noncrossing structures (Jin *et al.*, 2008; Huang *et al.*, 2009) has infinitely many shadows for any fixed  $k \geq 3$  (Reidys and Wang, 2010), and hence, like R&E, contains  $\gamma$ -structure with arbitrarily large  $\gamma$ . We note that every 1-structure is 4-noncrossing; more precisely, shadows (H) and (K) are 3-noncrossing, while shadow (L) and (M) consist of 3 mutually crossing arcs. See Fig. 7.

### 2.3 Minimum free energy folding of $\gamma$ -structures

We have shown in the previous section that 0-structures are simple RNA secondary structures. Their minimum free energy (MFE) configuration can be obtained by DP recursions (Waterman, 1978; Zuker and Stiegler, 1981) derived from a decomposition into suitable substructures. This decomposition can be expressed in terms of a context-free grammar (Dowell and Eddy, 2004; Steffen and Giegerich, 2005). In the simplest case, which corresponds to evaluating base pairs only, we consider a single non-terminal symbol  $S$



**Fig. 8.** Fragment-pairs in RNA structures: the rule  $I \rightarrow IA_1IB_1IA_2IB_2S$  induces the fragment-pairs  $[i_1, r_1]$ ,  $[s_1, j_1]$  and  $[i_2, r_2]$ ,  $[s_2, j_2]$ . Arcs connecting the two fragments of a pair are non-crossing, while arcs with both endpoints within the same fragment may be crossing such as those within  $[s_2, j_2]$ .

representing an arbitrary diagram over a segment and three terminal symbols to represent isolated vertices (symbol  $:$ ), openings (symbol  $($ ) and closings (symbol  $)$ ) of base pairs. We only need the three production-rules

$$S \rightarrow S, \quad S \rightarrow (S)S, \quad S \rightarrow \varepsilon, \quad (2.7)$$

to generate the corresponding language  $\mathcal{S}$ .

We shall use that (1) any 1-structure can be inductively generated from genus one structures and (2) that every genus one structure has shadow (H), (K), (L), or (M), to specify a multiple context-free grammar (MCFG) (Seki *et al.*, 1991). In contrast to context-free grammars, the non-terminal symbols of MCFGs may consist of multiple components which must be expanded<sup>3</sup> in parallel. In this way, it becomes possible to couple separated parts of a derivation and thus to generate crossings. In the case of 1-structures, the language  $\mathcal{S}$  is built upon sequences of intervals (*fragment-pairs*)  $[i, r]$ ,  $[s, j]$ , where  $(i, j)$ ,  $(r, s)$  are nested arcs. Arcs having endpoints in the different fragments are assumed to be noncrossing; see Fig. 8. For the MCFG, the fragments of a pair are associated with two different (coupled) components of a 2-dimensional non-terminal symbol.

Accordingly, we (re)introduce the following symbols:

- non-terminal  $S$ , representing secondary structure elements (i.e., diagrams without crossing arcs) according to the rules given above,
- non-terminals  $I$  and  $T$ , representing an arbitrary 1-structure,
- non-terminals  $\vec{X} = [X_1, X_2]$  with two components used to represent a fragment-pair with nested arcs,  $X \in \{H, K, L, M\}$ ,
- terminals  $(x, )_X$  denoting the opening and closing of a base pair, resp., where  $X$  is one of the types  $H, K, L$  or  $M$ .

Different brackets as well as the different non-terminals of pattern  $\vec{X}$  are used to distinguish nestings of the various kinds of shadows. Finally, we specify the production-rules of our unambiguous MCFG

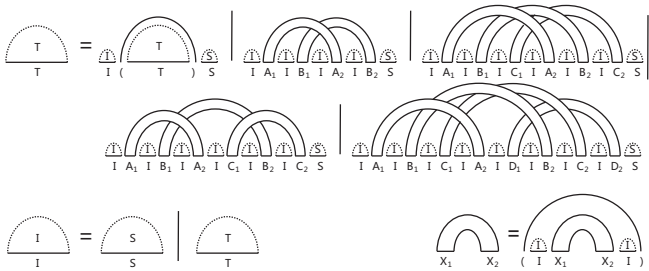
<sup>3</sup> This coupling is only required for components that were generated by the same production step. Components, even if of the same kind, derived in different steps are independent of each other.

$\mathcal{R}_1$ :

$$\begin{aligned}
 I &\rightarrow S \mid T \\
 S &\rightarrow (S)S \mid :S \mid \epsilon \\
 T &\rightarrow I(T)S \\
 T &\rightarrow IA_1IB_1IA_2IB_2S \\
 T &\rightarrow IA_1IB_1IA_2IC_1IB_2IC_2S \\
 T &\rightarrow IA_1IB_1IC_1IA_2IB_2IC_2S \\
 T &\rightarrow IA_1IB_1IC_1IA_2ID_1IB_2IC_2ID_2S \\
 \vec{X} &\rightarrow [(XIX_1, X_2I)_X] \mid [(X, )_X],
 \end{aligned}$$

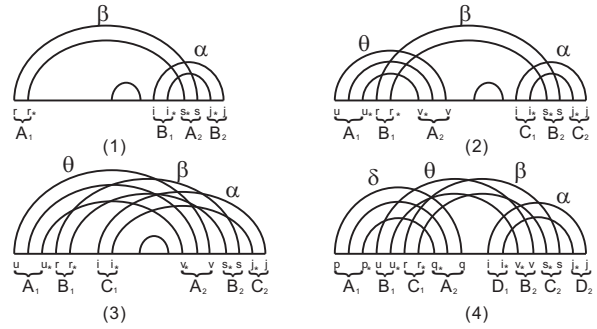
where  $X \in \{H, K, L, M\}$  distinguishes the four types of pseudo-knots.

**THEOREM 2.4.** Any RNA 1-structure can be **uniquely** decomposed via  $\mathcal{R}_1$ , and any diagram generated via  $\mathcal{R}_1$  is a 1-structure, see Fig. 9.



**Fig. 9.** Illustration of the grammar  $\mathcal{R}_1$ .

**PROOF.** We proceed by induction on the number of shadows. *Induction basis:* In a 1-structure  $\mathfrak{S}$  that contains no genus 1-shadow there are no crossings and hence the structure can be decomposed uniquely via the context-free grammar of secondary structures. *Induction step:* Suppose we are given a 1-structure containing  $r \geq 1$  shadows of genus one. We decompose from right to left. Everything is clear until we encounter a substructure containing a genus 1 shadow. For an arc  $\alpha = (i, j)$ , we distinguish two cases: (I)  $\alpha$  is not crossed, or (II)  $\alpha$  is crossed by another arc. In case of (I), there exists a 1-structure nested in  $\alpha$ . In case of (II), we consider the partial order  $\leq$ , where  $(i, j) \leq (r, s)$  if and only if  $r < i$  and  $j < s$ . Since crossing arcs in a 1-structure are contained in one of the four base types, we distinguish the following scenarios  
(H): then there exist maximal base pairs  $\beta = (r, s)$ , where  $r < i < s < j$ ,  
(K): then there exist maximal base pairs  $\beta = (r, s)$  and  $\theta = (u, v)$ , where  $u < r < v < i < s < j$ ,  
(L): then there exist maximal base pairs  $\beta = (r, s)$  and  $\theta = (u, v)$ , where  $u < r < i < v < s < j$ ,  
(M): then there exist maximal base pairs  $\beta = (r, s)$ ,  $\theta = (u, v)$  and  $\delta = (p, q)$ , where  $p < u < r < q < i < v < s < j$ .

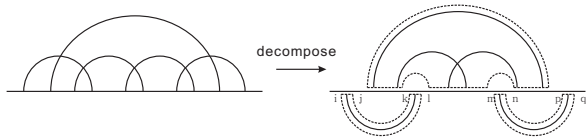


**Fig. 10.** Fragmentation: the four cases corresponding to the four shadows (H), (K), (L) and (M). In (1), there are two maximal arcs:  $\alpha = (i, j)$  and  $\beta = (r, s)$ , where  $r < i < s < j$ , whence the diagram has shadow (H). Here,  $\alpha_* = (i_*, j_*)$  is the minimal arc crossing  $C(\alpha)$  and  $\beta_* = (r_*, s_*)$  is the minimal arc crossing  $C(\beta)$ . We have  $B_1 = [i, i_*]$ ,  $B_2 = [j, j_*]$ ,  $A_1 = [r, r_*]$ ,  $A_2 = [s, s_*]$ . Cases (2), (3) and (4) are analyzed similarly.

Consider the set  $C(\alpha)$  of arcs that are crossed by  $\alpha$  and the minimal arc  $\alpha_*$  that crosses any element of  $C(\alpha)$ . Here minimality is considered with respect to the partial order  $\leq$ , where  $(i, j) \leq (r, s)$  if and only if  $r < i$  and  $j < s$ . It follows that  $\alpha = (i, j)$  and  $\alpha_* = (i_*, j_*)$  induce the fragment pair  $[i, i_*]$  and  $[j_*, j]$ . We similarly obtain the corresponding arcs  $\beta_*$ ,  $\theta_*$  or  $\delta_*$ , which induce at most four fragment-pairs and correspond to a unique shadow of type (H), (K), (L) or (M). See Fig. 10. By construction, the number of genus 1 shadows of any substructure contained in such a fragment-pair is reduced at least by one and can by induction hypothesis, be uniquely decomposes via  $\mathcal{R}_1$ . Finally, any structure generated via  $\mathcal{R}_1$  is constructed from top-to-bottom by iteratively building configurations of arcs having shadow (H), (K), (L) or (M). Thus any structure obtained via  $\mathcal{R}_1$  is indeed a 1-structure completing the proof of the theorem.

*2-structures.* A folding algorithm for 2-structures requires an analogous enumeration of all (irreducible) shadows of genus 2. From eq. (2.6), it is straightforward to explicitly derive the 21 shadows of genus 2 with 4 arcs, see SM Fig. 10. As in the case of genus 1, arc-insertions into these 21 configurations leads to the complete set of 3472 shadows of genus two. This large number makes it infeasible to build a practically useful folding algorithms for *all* 2-structures. It may be useful, however, to deal with a (small) subset of shadows. The complexity of such an algorithm is determined by the complexity of decomposing the individual shadows by means of MCFG-production rules reminiscent of those for  $\mathcal{R}_1$ . For instance, the shadow of the HDV structure displayed in Fig.6, (2), is contained in the R&E class and can therefore be computed in  $O(N^6)$  time and  $O(N^4)$  space. However, when resorting to our approach its time complexity is at least  $O(N^8)$ : the shadow presented in Fig. 11 requires a DP algorithm with  $O(N^8)$  time- and  $O(N^6)$  space-complexity. It is ongoing work to devise a sensible folding algorithm for 2-structures.

*MFE folding of 1-structures.* If we make use of a naïve table-based parsing scheme, checking for each subword  $s$  of the input and for each rule  $f$  whether  $f$  can produce  $s$ , a rule like  $f =$



**Fig. 11.** Folding of 2-structures: The shadow shown here is *not* contained in the R&E class of structures and cannot be generated by gap-matrices. It can be decomposed, however, using the 8 indexes  $i, j, k, l, m, n, p$  and  $q$ , thus implying a  $O(N^8)$  time-complexity. This makes use of a six-dimensional gap matrix  $G_{j,k,l,m,n,p}$ , which implies  $O(N^6)$  space-complexity.

$I \rightarrow IA_1IB_2IC_1IA_2ID_1IB_2IC_2ID_2S$  introduces a complexity  $O(N^{18})$ : First, we must process  $O(N^2)$  different subwords  $s$  induced by an input of size  $n$ . Second, each non-terminal but the first on the right-hand side of the production introduces an additional split point which specifies the part of  $s$  to be generated by the corresponding non-terminal. Since its location may freely be chosen within  $s$ , each split point gives rise to another loop variable, and hence contributes a factor  $O(N)$  to the runtime.

Even if there are much more sophisticated parsing algorithms, it is useful to consider this simple scheme since it directly translates into a recursion for a DP algorithm typically used to compute structures of minimum free energy. Furthermore, it is possible to introduce intermediate steps in the derivation of our language by making use of additional non-terminals and production-rules such that the time complexity can be reduced to  $O(N^6)$ . For that purpose let the non-terminal  $I'$  represent 1-structures in which no structures with shadow (H), (K), (L) or (M) are nested and the last vertex is paired. We introduce the non-terminal symbols  $\vec{U} = [U_1, U_2]$ ,  $\vec{V} = [V_1, V_2]$  and  $\vec{W} = [W_1, W_2]$  assumed to represent intermediate fragment-pairs and the production-rules

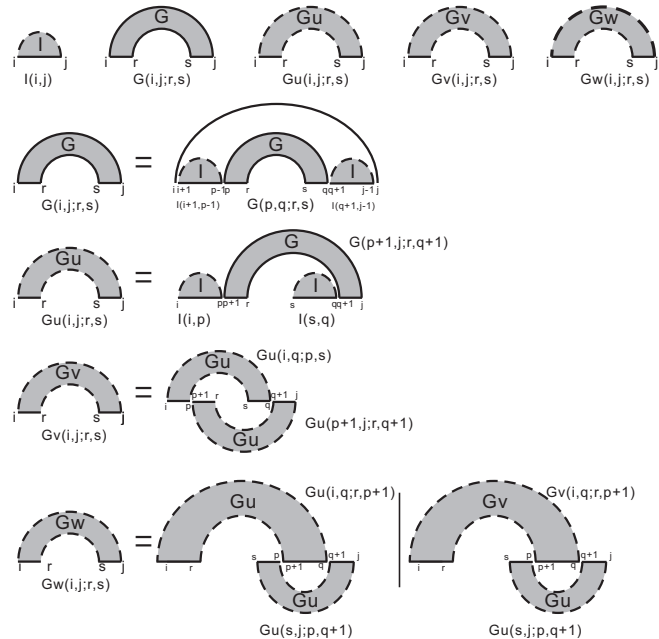
$$\begin{aligned} \vec{U} &\rightarrow [IX_1, IX_2] \\ \vec{V} &\rightarrow [U_1U'_1, U_2U'_2] \\ \vec{W} &\rightarrow [U_1, U'_1U_2U'_2] \mid [V_1, U_1V_2U_2] \end{aligned}$$

where  $(U'_1, U'_2)$  is a marked copy of  $(U_1, U_2)$  used to identify the components which must later be expanded in a coupled way. Accordingly, we replace the derivations of  $T$  in  $\mathcal{R}_1$  as follows:

$$\begin{aligned} T &\rightarrow I(T)S \mid I'S \\ I' &\rightarrow V_1V_2 \mid U_1V_1U_2V_2 \mid U_1W_1U_2W_2 \end{aligned}$$

Note that syntactically, i.e., considered as dot-bracket representations, the 1-structures can be generated by a MCFG, parsable in time  $O(N^5)$ . However, in that case, corresponding brackets are not generated in a coupled way making the grammar inappropriate for algorithmic purposes.

As typical for DP and in analogy to our parsing scheme, we use 2-dimensional matrices to store the optimal structure over a fragment. The matrix is indexed by the sequence coordinates of the endpoints. It can be a simple secondary structure  $S$  or a substructure of higher genus. For the fragment-pairs, i.e., for the non-terminals of dimension two, 4-dimensional matrices indexed by the endpoints of both linked fragments are required to store the optimal structure over them. Suppose the pair of fragments is  $[i, r]$  and



**Fig. 12.** The decomposition for 4-dimensional matrices  $G$ ,  $Gu$ ,  $Gv$ , and  $Gw$ .

$[s, j]$ , and let  $Gu(i, j; r, s)$  be the fragment-pair (associated with  $[U_1, U_2]$ ),  $Gv(i, j; r, s)$  be the fragment-pair  $[V_1, V_2]$ ,  $Gw(i, j; r, s)$  be the fragment-pair  $[W_1, W_2]$ , and  $G(i, j; r, s)$  be the fragment-pair  $[X_1, X_2]$ . The recursions for these matrices, summarized in graphical form in Fig. 12, are determined directly by the grammar.

We can conclude from the rewriting rules that the computation of the 2-dimensional matrices requires at most three loop variables, and there are  $O(N^2)$  many of them. Accordingly,  $O(N^5)$  operations are required to fill the associated 2-dimensional matrices. For the 4-dimensional matrices, two loop variables are needed for each of the corresponding rewriting rules (those with a left-hand side of dimension 2) for there are in each case two split points introduced by the right-hand sides of the corresponding productions. Since we need to compute  $O(N^4)$  matrix entries, the total run time is in  $O(N^6)$ . Obviously,  $O(N^4)$  space is required to store these tables. Accordingly, the algorithm can generate all 1-structures in  $O(N^6)$  time and  $O(N^4)$  space, i.e., with the same complexity as `pknotsRE` (Rivas and Eddy, 1999) (for the larger R&E class). The advantage of 1-structures is that structurally different shadows can be parametrized in different ways, and that the search space is restricted to moderately complex shadows. In contrast, the language of R&E-structures is based on crossings and can neither identify blocks of arcs not restrict the genus of the shadows. For more structure classes restricted to  $H$ -structures, `NUPACK` (Dirks and Pierce, 2003) requires  $O(N^5)$  time and  $O(N^4)$  space.

This is substantially more demanding, of course, than the  $O(N^4)$  time and  $O(N^2)$  memory complexity of `pknotsRG` Reeder and Giegerich (2004), which, however, deals with a very restricted subset of  $H$ -shadow structures, demanding that helices are maximally extended and perfect in the sense that they are not interrupted by bulge- or interior-loops. `pknotsRG` thus is not guaranteed to find



the minimum energy structure within the class H-shadow structures. A related fast heuristic treats the (K)-shadow as a superposition of the two H-shadows Theis *et al.* (2010).

## 2.4 Partition function and sampling

We have shown that the MCFG  $\mathcal{R}_1$  uniquely generates all 1-structures, i.e., it is unambiguous. Consequently,  $\mathcal{R}_1$  can be employed to count 1-structures over a given sequence  $x$  and to compute the corresponding partition function

$$Q = \sum_{s \in \mathfrak{S}_x} e^{-G(s)/RT},$$

where  $R$  is the universal gas constant,  $T$  is the temperature,  $G(s)$  is energy of structure  $s$  over sequence  $x$ , and  $\mathfrak{S}_x$  is the set of 1-structures in which all base pairs  $(i, j)$  satisfy the base pairing rules for RNA, i.e.,  $x_i x_j \in \{AU, UA, GC, CG, GU, UG\}$ . Let  $N_{i,j}$  denote the substructure represented by the nonterminal symbol  $N$  in  $\mathcal{R}_1$  over the fragment  $[i, j]$ , and let  $\vec{X}_{i,j;r,s}$  denote the fragment-pair  $\vec{X} = [X_1, X_2]$ , where  $X_1 = [i, r]$  and  $X_2 = [s, j]$  in the recursions for energy minimization. For each of these symbols, we introduce corresponding partial partition functions  $Q_{N_{i,j}}$  and  $Q_{\vec{X}_{i,j;r,s}}$ . Since the MCFG is unambiguous, the recursions for the partial partition functions are derived by replacing minima by sums and addition of energy contribution by multiplication of partial partition functions, see e.g., (Voß *et al.*, 2006). For instance, the recursion for the partition functions corresponding to the nonterminal symbol  $T$  reads

$$Q_{T_{i,j}} = \sum_h Q_{I'_{i,h}} \times Q_{S_{h+1,j}} + \sum_{h,\ell} Q_{I_{i,h-1}} \times Q_{T_{j+1,\ell-1}} \times Q_{S_{\ell+1,j}} \times e^{-E[h,\ell]/RT},$$

where  $E[h, \ell]$  denotes the energy of the loop closed by the base pair  $(h, \ell)$ .

The probabilities  $\mathbb{P}_{N_{i,j}}$  of partial structures of type  $N$  over the fragment  $[i, j]$  and the probabilities  $\mathbb{P}_{\vec{X}_{i,j;r,s}}$  of partial structures of type  $\vec{X}$  over the fragment pair  $[i, j], [r, s]$  are readily calculated from the partial partition functions. These “backward recursions” are analogous to those derived by McCaskill (1990) for crossing free structures: Let  $\Lambda_{N_{i,j}}$  be the set of 1-structures containing  $N_{i,j}$  and let  $\Lambda_{\vec{X}_{i,j;r,s}}$  be the set of 1-structures containing the fragment-pair  $\vec{X}_{i,j;r,s}$ . It follows that we have

$$\mathbb{P}_{N_{i,j}} = \sum_{s \in \Lambda_{N_{i,j}}} \mathbb{P}_s, \quad \mathbb{P}_{\vec{X}_{i,j;r,s}} = \sum_{s \in \Lambda_{\vec{X}_{i,j;r,s}}} \mathbb{P}_s.$$

Suppose  $N_{i,j}$  or  $\vec{X}_{i,j;r,s}$  are obtained by decomposing  $\theta_s$ . The conditional probabilities  $\mathbb{P}_{N_{i,j}|\theta_s}$  and  $\mathbb{P}_{\vec{X}_{i,j;r,s}|\theta_s}$  are then given by  $Q_{\theta_s}(N_{i,j})/Q_{\theta_s}$  and  $Q_{\theta_s}(\vec{X}_{i,j;r,s})/Q_{\theta_s}$  respectively. Here  $Q_{\theta_s}$  represents the partition function of  $\theta_s$ , and  $Q_{\theta_s}(N_{i,j})$  and  $Q_{\theta_s}(\vec{X}_{i,j;r,s})$  represent the partition functions for those  $\theta_s$ -configurations that contain  $N_{i,j}$  and  $\vec{X}_{i,j;r,s}$  respectively. Taking the sum over all possible  $\theta_s$ , we obtain

$$\mathbb{P}_{N_{i,j}} = \mathbb{P}_{\theta_s} \frac{Q_{\theta_s}(N_{i,j})}{Q_{\theta_s}}, \quad \mathbb{P}_{\vec{X}_{i,j;r,s}} = \mathbb{P}_{\theta_s} \frac{Q_{\theta_s}(\vec{X}_{i,j;r,s})}{Q_{\theta_s}}.$$

From this backward recursion, one immediately derives a stochastic backtracing recursion from the probabilities of partial structures that

generates a Boltzmann sample of 1-structures, see (Tacker *et al.*, 1996; Ding and Lawrence, 2003; Huang *et al.*, 2010) for analogous constructions.

The basic data structure for this sampling is a stack  $A$  which stores blocks of the form  $(i, j, N)$  (or  $(i, j; r, s, \vec{X})$ ), presenting substructures of nonterminal symbols  $N$  over  $[i, j]$  (or  $\vec{X}$  over  $[X_1, X_2]$  where  $X_1 = [i, r]$  and  $X_2 = [s, j]$ ).  $L$  is a set of base pairs storing those removed by the decomposition step in the grammar. We initialize with the block  $(1, n, I)$  in  $A$ , and  $L = \emptyset$ . In each step, we pick up one element in  $A$  and decompose it via the grammar with probability  $Q^M/Q^N$ , where  $Q^N$  is the partition function of the block which is picked up from  $A$ , and  $Q^M$  is the partition function of the target block which is decomposed by the rewriting rule. The base pairs which are removed in the decomposition step are moved to  $L$ . For instance, according to the rewriting rule  $T \rightarrow I(T)S$ , the block  $(i, j, T)$  is decomposed into the three blocks:  $(i, h-1, I)$ ,  $(h+1, \ell-1, T)$ ,  $(\ell+1, j, S)$  and one base pair  $(h, \ell)$  which is to be removed. For fixed indices  $h, \ell$ , where  $i \leq h < \ell \leq j$ , the probability of decomposing  $(i, j, T)$  reads

$$\mathbb{P}_{h,\ell} = \frac{Q_{I_{i,h-1}} \times Q_{T_{j+1,\ell-1}} \times Q_{S_{\ell+1,j}} \times e^{-E[h,\ell]/RT}}{Q_{T_{i,j}}}.$$

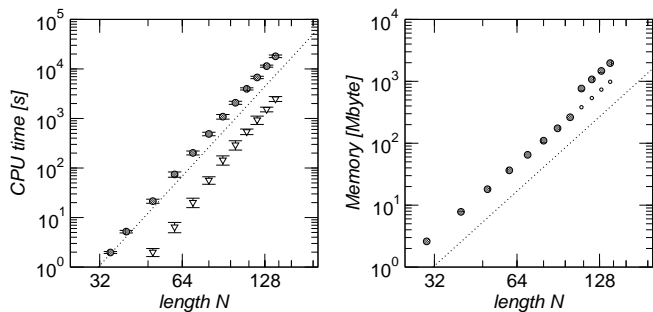
The sampling step is iterated until  $A$  is empty. The resulting 1-structure is the given by the list  $L$  of base pairs.

## 2.5 Software

*Implementation.* MFE folding, partition function including a computation of base pairing probabilities, and stochastic backtracing are implemented in `gfold`. The program is written in C.

*Energy Model.* Although the presentation above uses a simplified grammar that does not explicitly distinguish the usual loop types, `gfold` implements the Mathews-Turner energy model without dangles (Mathews *et al.*, 1999, 2004) for secondary structure elements. For pseudoknots, we use here an extended version of the Dirks-Pierce (DP) model (Dirks and Pierce, 2003) that allows different penalties  $\beta_X$  for the four topologically distinct pseudoknot types  $X = H, K, L, M$ . We have observed that the values of  $\beta_X$  have a substantial influence on the accuracy of the predicted structures. In both NUPACK and `pknotsRE`, a common pseudoknot penalty  $\beta_1$  is assigned whenever two gap matrices cross. Since the number of such crossings depends on the type of the pseudoknot, this algorithmic design would imply  $\beta_A = \beta_1$ ,  $\beta_B = \beta_C = 2\beta_1$ , and  $\beta_D = 3\beta_1$ . In `gfold`, these parameters are independent and can be adjusted to improve the performance. Since most experimentally known pseudoknots are of types (H) and (K), we focused in particular on the ratio of  $\beta_A$  and  $\beta_B$  and found that both sensitivity (the ratio of correctly predicted base pairs to the total number of base pairs in the reference structure) and positive predictive value reach a maximum for  $\beta_B = 1.3\beta_A$ . The pseudoknot penalty of type (H) coincides with that of the DP model, i.e.,  $\beta_A = \beta_1 = 9.6$  [kcal/mol]. The other penalties are set to  $\beta_B = 12.6$ ,  $\beta_C = 14.6$ , and  $\beta_D = 17.6$ ; see SM for details. An alternative set of pseudoknot parameters described by Andronescu *et al.* (2010) can easily be incorporated but would require a re-adjustment of these four topological penalties.

**Performance.** The current implementation of `gfold` is applicable to sequences with a length up to  $N \approx 150$  nucleotides on current PC hardware. Fig. 13 summarizes the resource requirements.



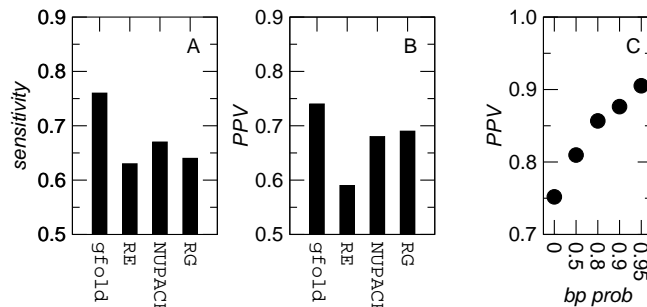
**Fig. 13.** Run time (left) and peak memory (right) of `gfold`. Timing information is given for MFE-only (triangles) and partition function with sampling 10,000 structures from the Boltzmann ensemble. To compute error bars, we folded between 10 ( $N > 100$ ) and 100 ( $N < 70$ ) randomly generated sequences on a Xeon E5410, 2.33Ghz, 48Gb memory. Memory allocation is independent of the sequence. For  $N \geq 100$ , double precision floats are necessary to avoid overflows. This leads to the jump in memory consumption by a factor of 2. Dotted lines indicate the theoretical behavior of  $O(N^6)$  (time) and  $O(N^4)$  (space). The slope for CPU time is slightly steeper than the theory since constraints among the 6 indices introduced by the minimum size of the complex pseudoknot elements lead to an additional speedup for small  $N$ .

We have observed that `gfold` provides a substantial increase in both sensitivity and a positive predictive value (PPV, ratio of correctly predicted base pairs to the total number of base pairs in the predicted structure) compared to the alternative DP approaches `pknotsRE` (Rivas and Eddy, 1999), `NUPACK` (Dirks and Pierce, 2003), and `pknotsRG-mfe` (Reeder and Giegerich, 2004), and that `gfold` provides a substantial increase in accuracy, cf. Fig. 14. In an evaluation on the entire `Pseudobase` (van Batenburg *et al.*, 2001), `gfold` achieves a sensitivity of 0.762 and PPV of 0.761. As detailed in SM (Tab.S-3), the performance varies substantially between different classes of sequences however. Interestingly, the more complex pseudoknots of type (K) are predicted with even higher accuracy (sensitivity 0.889, PPV 0.899) than the simpler, much more frequent type H.

The PPV of `gfold` predictions can be increased by filtering the base pairs of the MFE structure by their probability  $p$  of formation, which is computed by the partition function version of `gfold`. Accepting only base pairs with a predicted base pairing probability  $p > 0.95$  increases the PPV from 0.76 to more than 0.9, see Fig. 14C. In order to evaluate the false positive rate, we folded 100 tRNA sequences from Sprinzl’s tRNA database (Jühling *et al.*, 2009). `gfold` correctly identifies 94% of them as pseudoknot-free. In comparison, `NUPACK` correctly identifies 86% and `pknotsRG-mfe` 89% of this sample set.

### 3 DISCUSSION

Combinatorial models of pseudoknotted RNA structures are limited in two ways: On the one hand, exact algorithmic folding can



**Fig. 14.** Performance of `gfold`. Comparison of the average sensitivity (A) and PPV (B) of different prediction algorithms on a sample of 32 structures from `Pseudobase`. All details of this sample are given in the SM (Tab.S-2). (C) The PPV increases significantly if only base pairs with larger pairing probabilities as predicted by the partition function version of `gfold` are included in the predicted structure.

be constructed only for certain types of structures. On the other hand, the larger the structure sets are, the more base pairing patterns are contained in them that cannot be realized in nature due to steric constraints. Algorithm design so far has been mostly driven by the desire to reduce computational complexity. The idea behind `gfold`, in contrast, is to define a more suitable class of structures that can be generated by nesting and concatenating a small number of elementary building blocks. This recursive structure is captured by a fairly simple unambiguous multiple context-free grammar that translates in a canonical way to DP algorithms for computing the minimum energy structure and the partition function in  $O(N^6)$  time and  $O(N^4)$  space. In addition to MFE folding, we have implemented the computation of base pairing probabilities and a stochastic backtracing recursion, thus providing the major functionalities of RNA secondary structure prediction software for a very natural class of pseudoknotted structures.

The 1-structures considered here strike a balance between the generality necessary to cover almost all known pseudoknotted structures, and the restriction to topologically elementary structures that have a good chance to actually correspond to a feasible spatial structure. From a mathematical point of view, the characterization of structures in terms of irreducible components with given topological genus appears particularly natural and promises to reflect closely the ease with which a structure can be embedded in three dimensions. In addition, the grammar underlying `gfold` naturally distinguishes different types of pseudoknots and admits different energy parameters for them. We observe that this additional freedom of the parametrization leads to a substantial increase of sensitivity of type (K) pseudoknots, ( $0.63 \rightarrow 0.889$ ) and PPV ( $0.73 \rightarrow 0.899$ ) compared to the usage of a common penalty for each crossing of gap matrices. In terms of prediction accuracy, `gfold` thus compares favorably also with the leading alternative DP approaches to pseudoknotted structures.

**Acknowledgements.** This work was supported by the 973 Project of the Ministry of Science and Technology, the PCSIRT Project of the Ministry of Education, and the National Science Foundation of China to CMR and his lab, as well as the *Deutsche Forschungsgemeinschaft*, projects STA 850/2-1 & STA 850/7-1, the European Union FP-7 project QUANTOMICS (no. 222664) to PFS and his lab. JEA and RCP are supported by QGM, the Centre

for Quantum Geometry of Moduli Spaces, funded by the Danish National Research Foundation

## REFERENCES

- Akutsu, T. (2000). Dynamic programming algorithms for RNA secondary structure prediction with pseudoknots. *Discr. Appl. Math.*, **104**, 45–62.
- Andersen, J. E., Penner, R. C., Reidys, C. M., and Waterman, M. S. (2010). Enumeration of linear chord diagrams. *Comm. Pure and Appl. Math.* submitted.
- Andronescu, M. S., Pop, C., and Condon, A. E. (2010). Improved free energy parameters for RNA pseudoknotted secondary structure prediction. *RNA*, **16**, 26–42.
- Bailor, M. H., Sun, X., and Al-Hashimi, H. M. (2010). Topology links RNA secondary structure with global conformation, dynamics, and adaptation. *Science*, **327**, 202–206.
- Bon, M., Vernizzi, G., Orland, H., and Zee, A. (2008). Topological classification of RNA structures. *J. Mol. Biol.*, **379**, 900–911.
- Cai, L., Malmberg, R. L., and Wu, Y. (2003). Stochastic modeling of RNA pseudoknotted structures: a grammatical approach. *Bioinformatics*, **19**, 166–173.
- Chen, H.-L., Condon, A., and Jabbari, H. (2009). An  $O(n^5)$  algorithm for MFE prediction of kissing hairpins and 4-chains in nucleic acids. *J. Comp. Biol.*, **16**, 803–815.
- Chen, S. J. (2008). RNA folding: conformational statistics, folding kinetics, and ion electrostatics. *Annu Rev Biophys*, **37**, 197–214.
- Condon, A., Davy, B., Rastegari, B., Zhao, S., and Tarran, F. (2004). Classifying RNA pseudoknotted structures. *Theor. Comp. Sci.*, **320**, 35–50.
- Deogun, J. S., Donis, R., Komina, O., and Ma, F. (2004). RNA secondary structure prediction with simple pseudoknots. In *Proceedings of the second conference on Asia-Pacific bioinformatics (APBC 2004)*, pages 239–246. Australian Computer Society.
- Ding, Y. and Lawrence, C. E. (2003). A statistical sampling algorithm for RNA secondary structure prediction. *Nucleic Acids Res.*, **31**, 7280–7301.
- Dirks, R. M. and Pierce, N. A. (2003). A partition function algorithm for nucleic acid secondary structure including pseudoknots. *J. Comput. Chem.*, **24**, 1664–1677.
- Doudna, J. A. and Cech, T. R. (2002). The chemical repertoire of natural ribozymes. *Nature*, **418**, 222–228.
- Dowell, R. D. and Eddy, S. R. (2004). Evaluation of several lightweight stochastic context-free grammars for RNA secondary structure prediction. *BMC Bioinformatics*, **5**, 71.
- Ferré-D'Amaré, A. R., Zhou, K., and Doudna, J. A. (1998). Crystal structure of a hepatitis delta virus ribozyme. *Nature*, **395**, 567–574.
- Giedroc, D. P. and Cornish, P. V. (2009). Frameshifting RNA pseudoknots: structure and mechanism. *Virus Res.*, **139**, 193–208.
- Giegerich, R., Voß, B., and Rehmsmeier, M. (2004). Abstract shapes of RNA. *Nucl. Acids Res.*, **32**, 4843–4851.
- Haslinger, C. and Stadler, P. F. (1999). RNA structures with pseudo-knots: Graph-theoretical and combinatorial properties. *Bull. Math. Biol.*, **61**, 437–467.
- Hofacker, I. L., Fontana, W., Stadler, P. F., Bonhoeffer, L. S., Tacker, M., and Schuster, P. (1994). Fast folding and comparison of RNA secondary structures. *Monatsh. Chem.*, **125**, 167–188.
- Huang, F. W., Peng, W. W. J., and Reidys, C. M. (2009). Folding 3-noncrossing RNA pseudoknot structures. *J. Comp. Biol.*, **16**, 1549–1575.
- Huang, F. W. D., Qin, J., Reidys, C. M., and Stadler, P. F. (2010). Target prediction and a statistical sampling algorithm for RNA-RNA interaction. *Bioinformatics*, **26**, 175–181.
- Jin, E. Y., Qin, J., and Reidys, C. M. (2008). Combinatorics of RNA structures with pseudoknots. *Bull. Math. Biol.*, **70**, 45–67.
- Jühling, F., Mörl, M., Hartmann, R. K., Sprinzl, M., Stadler, P. F., and Pütz, J. (2009). tRNAdb 2009: compilation of tRNA sequences and tRNA genes. *Nucleic Acids Res.*, **37**, D159–D162.
- Kato, Y., Seki, H., and Kasami, T. (2006). RNA pseudoknotted structure prediction using stochastic multiple context-free grammar. *IPSJ Digital Courier*, **2**, 655–664.
- Kleitman, D. (1970). Proportions of irreducible diagrams. *Studies in Appl. Math.*, **49**, 297–299.
- Li, H. and Zhu, D. (2005). A new pseudoknots folding algorithm for RNA structure prediction. In L. Wang, editor, *COCOON 2005*, volume 3595, pages 94–103, Berlin. Springer.
- Loehl, M. and Moffatt, I. (2008). The chromatic polynomial of fatgraphs and its categorification. *Adv. Math.*, **217**, 1558–1587.
- Lyngsø, R. B. and Pedersen, C. N. (2000). RNA pseudoknot prediction in energy-based models. *J. Comp. Biol.*, **7**, 409–427.
- Massey, W. S. (1967). *Algebraic Topology: An Introduction*. Springer-Verlag, New York.
- Mathews, D., Sabina, J., Zuker, M., and Turner, D. (1999). Expanded sequence dependence of thermodynamic parameters improves prediction of RNA secondary structure. *J. Mol. Biol.*, **288**, 911–940.
- Mathews, D., Disney, M., Childs, J., Schroeder, S., Zuker, M., and Turner, D. (2004). Incorporating chemical modification constraints into a dynamic programming algorithm for prediction of RNA secondary structure. *Proc Natl Acad Sci*, **101**, 7287–7292.
- Matsui, H., Sato, K., and Sakakibara, Y. (2005). Pair stochastic tree adjoining grammars for aligning and predicting pseudoknot RNA structures. *Bioinformatics*, **21**, 2611–2617.
- McCaskill, J. S. (1990). The equilibrium partition function and base pair binding probabilities for RNA secondary structure. *Biopolymers*, **29**, 1105–1119.
- Metzler, D. and Nebel, M. E. (2008). Predicting RNA secondary structures with pseudoknots by MCMC sampling. *J. Math. Biol.*, **56**, 161–181.
- Namy, O., Moran, S. J., Stuart, D. I., Gilbert, R. J. C., and Brierley, I. (2006). A mechanical explanation of RNA pseudoknot function in programmed ribosomal frameshifting. *Nature*, **441**, 244–247.
- Nebel, M. E. and Weinberg, F. (2011). An algebraic approach to RNA pseudoknotted structures. *submitted*.
- Nussinov, R., Piecznik, G., Griggs, J. R., and Kleitman, D. J. (1978). Algorithms for loop matching. *SIAM J. Appl. Math.*, **35**(1), 68–82.
- Penner, R. C., Knudsen, M., Wiuf, C., and Andersen, J. E. (2010). Fatgraph models of proteins. *Comm. Pure Appl. Math.*, **63**, 1249–1297.
- Reeder, J. and Giegerich, R. (2004). Design, implementation and evaluation of a practical pseudoknot folding algorithm based on thermodynamics. *BMC Bioinformatics*, **5**, 104.
- Reidys, C. M. and Wang, R. (2010). Shapes of RNA pseudoknot structures. *J. Comput. Biol.*, **17**, 1575–1590.
- Rivas, E. and Eddy, S. R. (1999). A dynamic programming algorithm for RNA structure prediction including pseudoknots. *J. Mol. Biol.*, **285**, 2053–2068.
- Rivas, E. and Eddy, S. R. (2000). The language of RNA: A formal grammar that includes pseudoknots. *Bioinformatics*, **16**, 334–340.
- Rødland, E. A. (2006). Pseudoknots in RNA secondary structures: Representation, enumeration, and prevalence. *J. Comp. Biol.*, **13**, 1197–1213.
- Seki, H., Matsumura, T., Fujii, M., and Kasami, T. (1991). On multiple context free grammars. *Theor. Comp. Sci.*, **88**, 191–229.
- Staple, D. W. and Butcher, S. E. (2005). Pseudoknots: RNA structures with diverse functions. *PLoS Biol.*, **3**, e213.
- Steffen, P. and Giegerich, R. (2005). Versatile and declarative dynamic programming using pair algebras. *BMC Bioinformatics*, **6**, 224.
- Tabaska, J. E., Cary, R. B., Gabow, H. N., and Stormo, G. D. (1998). An RNA folding method capable of identifying pseudoknots and base triples. *Bioinformatics*, **14**, 691–699.
- Tacker, M., Stadler, P. F., Bomberg-Bauer, E. G., Hofacker, I. L., and Schuster, P. (1996). Algorithm independent properties of RNA structure prediction. *Eur. Biophys. J.*, **25**, 115–130.
- Taufel, M., Licon, A., Araiza, R., Mireles, D., van Batenburg, F. H. D., Gulyaev, A., and Leung, M.-Y. (2009). PseudoBase++: an extension of PseudoBase for easy searching, formatting and visualization of pseudoknots. *Nucleic Acids Res.*, **37**, D127–D135.
- Theimer, C. A., Blois, C. A., and Feigon, J. (2005). Structure of the human telomerase RNA pseudoknot reveals conserved tertiary interactions essential for function. *Mol. Cell*, **17**, 671–682.
- Theis, C., Janssen, S., and Giegerich, R. (2010). Prediction of RNA secondary structure including kissing hairpin motifs. *Algorithms in Bioinformatics*, **6293**, 52–64.
- Uemura Y., Hasegawa, A., Kobayashi, S., and Yokomori, T. (1999). Tree adjoining grammars for RNA structure prediction. *Theor. Comp. Sci.*, **210**, 277–303.
- van Batenburg, F. H. D., Gulyaev, A. P., and Pleij, C. W. A. (2001). PseudoBase: structural information on RNA pseudoknots. *Nucleic Acids Res.*, **29**, 194–195.
- Vernizzi, G. and Orland, H. (2005). Large- $N$  random matrices for RNA folding. *Acta Phys. Polon.*, **36**, 2821–2827.
- Voß, B., Giegerich, R., and Rehmsmeier, M. (2006). Complete probabilistic analysis of RNA shapes. *BMC Biology*, **4**, 5.
- Waterman, M. S. (1978). Secondary structure of single-stranded nucleic acids. *Adv. Math. (Suppl. Studies)*, **1**, 167–212.
- Zuker, M. (1989). On finding all suboptimal foldings of an RNA molecule. *Science*, **244**, 48–52.
- Zuker, M. and Stiegler, P. (1981). Optimal computer folding of larger RNA sequences using thermodynamics and auxiliary information. *Nucleic Acids Res.*, **9**, 133–148.