Folding 3-noncrossing RNA pseudoknot structures

Fenix W.D. Huang*, Wade W.J. Peng* and Christian M. Reidys*,

 \star Center for Combinatorics, LPMC-TJKLC
† College of Life Sciences

Nankai University, Tianjin 300071, P.R. China

Phone: *86-22-2350-6800, Fax: *86-22-2350-9272

June, 2009

duck@santafe.edu

Abstract

In this paper we present the novel *ab initio* folding algorithm **cross**, which generates minimum free energy (mfe), 3-noncrossing, canonical RNA structures. Here an RNA structure is 3-noncrossing if it does not contain three or more mutually crossing arcs and canonical, if each of its stacks has size greater or equal than two. Our notion of mfe-structure is based on a specific concept of pseudoknots and respective loop-based energy parameters. The algorithm decomposes into three subroutines: first the inductive construction of motifs and their associated shadows, second the generation of the (rooted) skeleta-trees and third the saturation of the skeleta via context dependent dynamic programming routines.

Keywords: RNA pseudoknot structure, k-noncrossing, motif, skeleton, dynamic programming.

1. Introduction

In this paper we introduce the *ab initio* folding algorithm cross. Cross folds RNA (ribonucleic acid) sequences [52] into minimum free energy (mfe), 3-noncrossing pseudoknot structures. We give a selfcontained analysis of the algorithm, whose source code is publicly available at

www.combinatorics.cn/cbpc/cross.html

RNA exhibits a variety of 3-dimensional structural configurations, the so called tertiary structures, determining the functionality of the molecule, see Fig.1. Besides the noncrossing base pairings found in RNA secondary structures there exist further, cross-serial nucleotide interactions [48, 6, 56]. These bonds are called pseudoknots and occur in functional RNA, like, for instance, RNAseP [33], as well as ribosomal RNA [32]. In fact, RNA exhibits a diversity of biochemical capabilities [3], proved by the discovery of catalytic RNAs, or ribozymes [33], in 1981. Like proteins, RNA is capable of catalyzing reactions whereas transfer RNA acts as a messenger between DNA and protein.

These RNA functionalities motivate the study of the RNA structure prediction problem. The first mfe-folding algorithms for RNA secondary structures are due to [14, 31, 10] and the first dynamic programming (DP) folding routines for secondary structures were given by Waterman et al. [49, 55, 57, 37], predicting the loop-based mfe-secondary structure [52] in $O(n^3)$ -time and $O(n^2)$ -space. The general problem of RNA structure prediction under the widely used thermodynamic model is known to be NP-complete when the structures considered include arbitrary pseudoknots [34]. There exist however, polynomial time folding algorithms, predicting specific types of pseudoknots: Rivas et al. [43], Uemura et al. [53], Akutsu [4] and Lyngsø [34]. There appears to be no consensus of what pseudoknot structures should be computed. As for the ab initio folding of pseudoknot RNA, we find the following two paradigms: Rivas and Eddy's [43] gap-matrix variant of Waterman's DP-folding routine for secondary structures [49, 54, 23, 55, 37], maximum weighted matching algorithms [13, 15] and the latter taylored for pseudoknot prediction [7, 50]. The former method folds into a somewhat "mysterious" class of pseudoknots [44] in polynomial time. However, Andronescu et al. [5] as well as Akutsu [4] have shown that output classes of DP-algorithms can be well-defined in terms of certain grammars. Further variants of DP-algorithms have been developed by Dirks and Pierce [11], Reeder and Giegerich [39] and Ren et al. [42]. Further ideas in the context of pseudoknot folding involve the iterated loop matching approach [45] and the sampling of RNA structures via the Markov-chain Monte-Carlo method [36].

Pseudoknot folding algorithms employing the DP-paradigm, see Section 6.1, inevitably produce arbitrarily high crossing numbers, see Tab.1 and Fig.2. Despite this, they cannot generate nonplanar 3-noncrossing RNA structures [38]. The generation of high crossing numbers is insofar problematic as it implies a very large output class. Already for k=4, i.e. for RNA structures exhibiting three mutually crossing arcs, we have an exponential growth rate of 6.8541–exceeding that of the number of natural sequences. This means, that only an exponentially small fraction of them is actually realized by some sequence. In this context, we remark that this growth rate appears to grow linearly in k, see Tab.1. As a result, any type of study, along the lines of [47, 26, 46, 41, 21, 17, 18], is purely computational and does not allow to draw conclusions in the sense of [51].

Let us next describe the output class of **cross**. To this end we introduce k-noncrossing diagrams, i.e. graphs over [1,n] with vertex degrees ≤ 1 , represented by drawing the vertices in a horizontal line and its arcs (i,j), where i < j, in the upper half-plane, containing at most k-1 mutually crossing arcs. The idea here is that the vertices and arcs correspond to nucleotides and Watson-Crick (**A-U**, **G-C**) and (**U-G**) base pairs, respectively. Diagrams have the following three key parameters: the maximum number of mutually crossing arcs, k-1, the minimum arc-length, λ and minimum stack-length, σ . The length of an arc (i,j) equals j-i and a stack of length σ is the sequence of "parallel" arcs of the form

$$((i,j),(i+1,j-1),\ldots,(i+(\sigma-1),j-(\sigma-1))),$$

see Fig.3. Cross generates RNA pseudoknot structures that are 3-noncrossing, $\sigma \geq 3$ -canonical, and that have a minimum arc-length $\lambda \geq 4$. The 3-noncrossing property is mostly for algorithmic convenience and the generalization to higher crossing numbers represents not a major obstacle. The class of k-noncrossing RNA structures is motivated by the observation that most natural RNA structures exhibit low crossing numbers [19]. The concept of k-noncrossing RNA structures is based on the combinatorial work of Chen et al. [8, 9]. Since we are interested in mfe-structures, we consider 3-canonical structures, i.e. those in which each stack has length at least three. The minimum arc-length of four is also a result of biophysical constraints. In Fig.4 we present a particular 3-noncrossing, 3-canonical RNA structure: the natural structure of the HDV-virus and the structure as folded by cross.

Theorem 1 implies exact enumeration results as well as an array of exponential growth rates indexed by k and σ . The latter are presented in Tab.2 and are of relevance in the context of the asymptotic analysis of the algorithm. Tab.2 furthermore shows that 3-noncrossing, σ -canonical RNA structures have remarkably moderate growth rates. σ -canonical structures with higher crossing numbers exhibit also moderate growth rates, indicating that generalizations of the current implementation of cross for 4- or 5-noncrossing canonical RNA structures are feasible.

Our energy model is a generalization of [43, 11]. We remark that we do not consider dangles, see Fig.5. It is indeed a generalization, since we consider 3-noncrossing nonplanar structures, see Section 6.4, which cannot be generated by the current

implementation of gap-matrices [38]. As for the pseudoknot energy parameter we have

$$G^{pseudo} = \beta_1 + B \cdot \beta_2 + U \cdot \beta_3,$$

where β_1 , β_2 and β_3 parameterize specific penalties, B is the number of base pairs and U is the number of unpaired bases therein. In Fig.6 we illustrate the energy parameter of a non-planar configuration.

By a priori generating 3-noncrossing RNA structures [27, 28, 35], cross has a transparent output class – a feature shared only by secondary structure folding algorithms. For k = 3, ..., 9 it is possible to prove via central limit theorems [29, 25] that, irrespective of energy parameters [40, 24], a variety of generic properties of sequence-structure maps into RNA pseudoknot structures hold. We can conclude from this, that cross generates exponentially many different 3-noncrossing structures having a neutral network of exponential size. Even in its current implementation, i.e. restricted to 3-noncrossing structures cross is capable of generating non-planar configurations. The extension of cross to k > 3 is work in progress.

In Fig.7 we display a nonplanar structure folded by **cross**. The mfe of this structure is given by -41.8 kcal/mol.

2. Outline

Cross is an *ab initio* folding algorithm that maps RNA sequences into 3-noncrossing RNA structures. It is guaranteed to search all 3-noncrossing, σ -canonical structures and derives some (not necessarily unique), loop-based mfe-configuration.

In Fig.8 we present an overview of cross. The algorithm calls three subroutines:

- Shadow: a procedure that generates all irreducible shadows for a given input sequence,
- SkeletonBranch: generating the skeleta-trees rooted in irreducible shadows and
- Saturation: a DP-filling of the skeleta.

The subroutines are described in detail in Sections 3, 4 and 5. Fig.8 provides an outline of these three subroutines. Before going into a detailed description, let us mention the key ideas and give the pseudocode of **cross**: here **Optimal** (S) is a procedure that selects the best (mfe) structure from S. Combine (A, B) combines structure A and B, which means to take the union of their arc-sets. IS, Skeleton, OS, OSM denote matrices which store substructures over regions [i, j]. Each entry in these matrices represents in fact a set of structures.

As for the general strategy, **cross** constructs 3-noncrossing RNA structure "from top to bottom".

Phase I (Shadow): Here for given sequence a set of motifs, representing somewhat maximal stacks of a structure are generated. Despite the fact that motifs can exhibit

complicated crossings, they are *inductively* generated, see Section 6.3. The motifs are then "shadowed", i.e. their stacks are extended "from top to bottom". The number of stacks is a key parameter of an irreducible shadow which allows for their inductive generation.

Phase II (SkeletonBranch): Given an irreducible shadow, the second step of cross consists in generating, the skeleta-tree. The nodes of this tree are particular 3-noncrossing structures, obtained by successive insertions of stacks, see Section 4. Intuitively, a skeleton encapsulates all cross-serial arcs that cannot be recursively computed. We describe the insertion scheme in Section 4 and prove via Proposition 1 and 2 that our data-structures are trees and capable of generating all skeleta. Here we control the tree complexity via the (total) number of pseudoknots.

Phase III (Saturation): In the third subroutine each skeleton is saturated via DP-routines. In Section 5 we describe in detail the context-dependent DP-routines via which specific substructures are being inserted. After the saturation of a skeleton we have arrived at the final mfe-3-noncrossing structure.

The design of cross is centered at the generation of four matrices: IS, Skeleton, OS and OSM. IS is generated by Shadow, and IS[i][i+j] stores all shadows whose origin and terminus are i and i+j, respectively. Skeleton is generated by SkeletonBranch, and Skeleton[i][i+j] stores the trees rooted from IS[i][i+j]. The matrices OS and OSM are both generated by Saturation, which ramifies quickly into

many subcases, depending on the loop-context encountered, see Section 6.5. Consequently cross is fundamentally different from the pseudoknot DP-routines found in the literature [43], see Section 6.1. The DP-routine of [43] can neither create nonplanar 3-noncrossing structures nor control the maximal number of mutually crossing arcs (crossing number). In fact, by construction, DP-routines are limited to generate pseudoknot complexity by increasing the crossing number. As nonplanar 3-noncrossing structures indicate, see Fig.7, structural complexity is not tantamount to the latter.

Shadow and SkeletonBranch have by design exponential time complexity. Only Saturation is polynomial time. Beyond the asymptotic analysis of motifs, given in Section 3, further analysis, in particular that of skeleta trees, is work in progress. In Fig.9 we display the logarithm of the mean folding times obtained by folding 1000 random sequences as a function of n. These data suggest folding times with exponential growth rates of ≈ 1.146 and ≈ 1.254 , for 3-canonical and 4-canonical structures, respectively. Furthermore, in Fig.10 we display the distribution of the logarithm of folding times of 800 random random sequences of length 100. As for comparisons of structures predicted by cross versus natural structures contained in PseudoBase [1], see Tab.3.

3. Motifs and shadows

In this section we introduce motifs and their shadows which form the foundation of the structures generated by **cross**. Let \prec denote the partial order over the set of arcs (written as (i, j), i < j) of a k-noncrossing diagram, given by

$$(3.1) (i_1, j_1) \prec (i_2, j_2) \iff i_2 < i_1 \land j_1 < j_2.$$

A k-noncrossing core is a k-noncrossing diagram without any two arcs of the form (i,j), (i+1,j-1), see Fig.11. Any k-noncrossing RNA structure, S has a unique k-noncrossing core, c(S) [28], obtained in two steps: first one identifies all arcs contained in stacks, inducing a contracted diagram and secondly one relabels the vertices. Note that the core-map does in general not preserve arc-length.

A motif, \mathfrak{m} , is a k-noncrossing, σ -canonical structure, such that: (M1) \mathfrak{m} has a nonnesting core and (M2) all \mathfrak{m} -arcs are contained in stacks of length exactly $\sigma \geq 3$ and arc-length $\lambda \geq 4$. The set of all motifs is denoted by $\mathbb{M}_k^{\sigma}(n)$ and we set $\mu_{k,\sigma}^*(n) = |\mathbb{M}_k^{\sigma}(n)|$. (M1) is obviously equivalent to: all arcs of the core, $c(\mathfrak{m})$, are \prec -maximal, see Fig.12

Corollary 1. Suppose k = 3 and $\sigma \geq 2$, then motifs can be inductively generated and we have the asymptotic formula

(3.2)
$$\mu_{3,\sigma}^*(n) \sim c_\sigma \left(\zeta_\sigma\right)^{-n},$$

where c_{σ} and ζ_{σ}^{-1} are given in Tab.4.

The above corollary is based on Proposition 3, Section 6.3 and provides a constructive method for obtaining the motifs over a given sequence. Furthermore its shows feasibility via the asymptotic analysis of the space of motifs.

Let S be a 3-noncrossing, σ -canonical structure. We wish to make precise what it means to extend the stacks of a motif "from top to bottom", see Fig.13. To this end let (i, j) be a minimal arc in a motif-stack. Extending the stack then consists in successively drawing the arcs $(i+1, j-1), (i+2, j-2), \ldots$ A shadow of a motif is a 3-noncrossing, σ -canonical structure, derived by extending some of the motif-stacks.

4. The skeleta-tree

In this section we enter the second phase of cross. What will happen here, is that each irreducible shadow, generated during the first phase described in Section 3, gives rise to a tree of skeleta. The intuition behind this construction is that each treevertex, i.e. each skeleton, represents a maximal "non-inductive" arc configuration. This does not mean that a skeleton contains all crossing arcs of the final structure, but all further crossings are derived by adding independent substructures. In other words: their energy contributions are additive.

Suppose we are given a structure Q. Let α be a Q-arc and denote the set of Q-arcs that cross α by $\mathscr{A}_Q(\alpha)$. A skeleton, S, is a 3-noncrossing structure whose core has no noncrossing arcs, i.e. for any arc α we have $\mathscr{A}_S(\alpha) \neq \emptyset$ and whose dependency graph is connected, see Fig.14. In addition, in a skeleton over the region [i,j], $S_{i,j}$, the positions i and j are paired. Recall that an interval is a sequence of consecutive, unpaired bases $(i,i+1,\cdots,j)$, where i-1 and j+1 are paired. Furthermore, recall that a stack of length σ (see eq. (1.1)) is a sequence of parallel arcs $((i,j),(i+1,j-1),\ldots,(i+(\sigma-1),j-(\sigma-1)))$, which we write as (i,j,σ) . Note that $\sigma \geq \sigma_0$, where σ_0 is the minimum stack-length of the structure, see Fig.14. An irreducible shadow over [i,j] is denoted by $\mathrm{IS}(i,j)$. It is a particular skeleton, i.e. a skeleton in which there are no nested stacks.

We are now in position to construct the skeleta-tree. Suppose we are given a 3-noncrossing skeleton, S. We label the S-intervals $\{I_1, \ldots, I_m\}$ from left to right and consider pairs (S, r), where r is an integer. Given a pair (S, r) we construct new pairs (S', r') where $r' \geq r$ as follows: we replace a pair of intervals (I_p, I_q) , $i \in I_p$, $j \in I_q$, $i \geq r$ by the stack $\alpha = (i, j, \sigma)$, subject to the following conditions

- S' is a 3-noncrossing skeleton
- $(i + \sigma 1, j \sigma + 1)$ is a minimal element in (S', \prec)
- r' is the label of the first paired base preceding the interval I_p .
- i-1 and j+1 are not paired to each other.

Fig.15 displays the two basic scenarios via which stacks are being inserted. We refer to the above procedure as (i, j, σ) -insertion, denoting it by

$$(4.1) (S,r) \Rightarrow_{(i,j,\sigma)} (S',r').$$

Given a pair (S, r), subsequent insertions induce a directed graph, $G_{(S,r)}$, whose vertices are pairs (S', r') and whose (directed) arcs are given by

(4.2)
$$((S,r),(S',r')), \text{ where } (S,r) \Rightarrow_{(i,j,\sigma)} (S',r').$$

Remark 1. Note that it is checked whether or not (i, j, σ) can be added, i.e. (1) the bases $\{i, i+1, \dots, i+\sigma-1, j-\sigma+1, \dots, j-1, j\}$ are indeed unpaired and (2) (i-1, j+1) is not a base pair. The second property guarantees that the core of the stack (i, j, σ) is an arc in the core of S'.

We proceed by showing that $G_{(S,r)}$ is in fact a tree. In other words, the insertion-procedure is an unambiguous grammar.

Proposition 1. Let $T_1 = \{S \mid \exists r; (S,r) \in T\}$ and S_0 be a 3-noncrossing skeleton.

- (a) $G_{(S_0,r_0)}$ is a tree and for any two different vertices (S'_1,r'_1) and (S'_2,r'_2) in $G_{(S_0,r_0)}$, we have $S'_1 \neq S'_2$.
- (b) For k > 3, the graph morphism $\pi : \mathbb{T} \longrightarrow \mathbb{T}_1$, given by $\pi((S, r)) = S$ is not bijective.

Remark 2. For any k > 3, $G_{(S_0,r_0)}$ is a tree. However assertion (b) indicates that it is *really* a tree of pairs. That means, stack-insertions will in general generate two different pairs with equal first coordinate.

Next we prove that our unambiguous grammar indeed generates all skeleta, see Fig.16.

Proposition 2. Suppose we are given an irreducible shadow $S_0 = IS(i, j)$. Let $\mathbb{T}(S_0) = G_{(S_0,i)}$ denote its skeleton-tree and let $\mathbb{S}(S_0)$ be the set of all skeleta, that contains S_0 and whose maximal arcs are contained in S_0 . Then we have

Remark 3. In cross we control the complexity of the skeleta-tree by limiting the number of pseudoknots. This number is an input parameter whose default is set to two.

5. SATURATION

In this section we discuss the third phase of cross. The skeleta-trees constructed in the second phase organized the non-inductive substructures of an irreducible shadow derived in phase one. The objective of Saturation is to inductively "fill"

the remaining intervals of a given skeleton with specific substructures. Basically, all routines employed here follow the DP-paradigm.

This DP-routine is formulated via two mfe-structures over a region [i, j], OSM(i, j) and OS(i, j). OSM(i, j) is the mfe-structure derived from a skeleton S over the region [i, j], and OS(i, j) is the mfe-structure over the region [i, j].

As for the construction of OS(i, j) via OSM(i', j'), see Fig.17, we consider position i in OS(i, j). If i is paired, then i is contained in some OSM(i, s). Then OS(i, j) induces a substructure S_2 over [s + 1, j]. By construction $OS(i, j) = OSM(i, s) \dot{\cup} S_2$, whence $S_2 = OS(s + 1, j)$ and in particular we have

(5.1)
$$\forall i < s < j, \quad \epsilon(OS(i,j)) = \epsilon(OSM(i,s)) + \epsilon(OS(s+1,j)),$$

where $\epsilon(S)$ denotes the energy of a 3-noncrossing structure S. Suppose next i is unpaired in OS(i, j). Since ϵ is a loop-based energy, we can conclude $OS(i, j) = \{\emptyset\} \dot{\cup} OS(i+1, j)$, i.e. we have

(5.2)
$$\epsilon(OS(i,j)) = \epsilon(OS(i+1,j)) + Q$$

where Q represents the energy contribution of a single, unpaired nucleotide. Accordingly, we can inductively construct OS(i, j) via the criterion

$$\forall i < s \le j, \quad \epsilon(\mathrm{OS}(i,j)) = \min\{\epsilon(\mathrm{OS}(i+1,j)) + Q, \epsilon(\mathrm{OSM}(i,s)) + \epsilon(\mathrm{OS}(s+1,j))\}.$$

Next we show how to obtain OSM(i, j) by the Saturation subroutine: for a given skeleton over [i, j], S(i, j), Saturation (S(i, j)) is a procedure that obtain OSM(i, j). Since S(i, j) is a skeleton, OSM(i, j) is irreducible. In order to obtain OSM(i, j), we simply insert the mfe-substructure in each interval of S(i, j).

The insertion of a substructure induces loops. Therefore, similar to the DP-routine for secondary structures, we have three types of mfe-substructure, $OS_{mul}(i,j)$, $OS_{pk}(i,j)$ and $OS_0(i,j)$. They satisfy identical recursions but subject to different energy parameters. Here $OS_{mul}(i,j)$ represents the optimal substructure nested in a multi-loop and $OS_{pk}(i,j)$ the one nested in a pseudoknot and $OS_0(i,j)$ the optimal substructure, otherwise. We present these recursions in Section 6.5.

Now we can inductively construct the arrays of structures OS(i, j) and OSM(i, j) via OS and OSM structures over smaller intervals. As a result, we finally obtain the structure OS(1, n), i.e. the mfe-structure.

6. Appendix

In this section we discuss the DP-routines in the context of pseudoknot RNA structures (Section 6.1), the combinatorics of k-noncrossing canonical RNA structures and motifs (Section 6.2 and Section 6.3), loops including the unique loop-decomposition (Section 6.4) and finally give all recursions of Saturation (Section 6.5).

6.1. **Pseudoknots via DP-routines.** In this section we discuss the main ideas behind the DP-paradigm in the context of pseudoknot folding by means of analyzing the algorithm of Rivas and Eddy [43, 44, 12]. The key observation here is the use of gap-matrices in addition to the wx and vx, discussed above, see Fig.18.

There are four gap-matrices, whx(i, j, r, s), vhx(i, j, r, s), yhx(i, j, r, s) and zhx(i, j, r, s), given in Tab.5.

The algorithm coincides with the DP-routine for secondary structures in case of gaps of size zero, that is r = s - 1. Then

(6.1)
$$whx(i, j; r, r+1) = wx(i, j)$$

(6.2)
$$zhx(i, j; r, r+1) = vx(i, j),$$

for $i \leq k \leq j$. In principle, any number of gap-matrices can be employed. However, the algorithm, in its current implementation, is truncated at O(whx + whx + whx), that is, at each step at most two gap-matrices are used. We present the basic recursions in Fig.19.

Via those recursions one computes the four gap-matrices whx(i, j, r, s), vhx(i, j, r, s), yhx(i, j, r, s) and zhx(i, j, r, s). Then, analogous to the case of secondary structures [55], we obtain a mfe-configuration via backtracking.

In Fig.20 we showcase a nonplanar 3-noncrossing structure, which cannot be generated by two gap-matrices [38].

While the inductive formation of two (or more) gap-matrices generates arbitrarily high numbers of mutually crossing arcs, see Fig.20, this method fails to generate non-planar, 3-noncrossing pseudoknots. In Fig.7, we give an example of a 3-noncrossing structure generated by cross, that cannot be constructed using two gap-matrices. It is clear, that gap-matrices can and will generate nonplanar arc configurations. However, they can only facilitate this via increasing the crossing number. Fig.7 makes evident that the situation is real and more complex: nonplanarity is not tied to crossings—there are planar as well as nonplanar 3-noncrossing structures. The situation becomes much more involved for higher crossing numbers.

6.2. Combinatorics of RNA pseudoknot structures. Let $\mathsf{T}_{k,\sigma}^{[4]}$ denote the number of k-noncrossing, σ -canonical RNA structures over [n] with minimum arc length four $(\langle k, \sigma \rangle$ -structure). The generating function [35],

$$\mathbf{T}_{k,\sigma}^{[4]}(z) = \sum_{n>0} \mathsf{T}_{k,\sigma}^{[4]}(n) z^n \quad k, \sigma \ge 3$$

is closely related to $\mathbf{F}_k(z) = \sum_n f_k(2n,0)z^{2n}$, the ordinary generating function of k-noncrossing matchings. Beyond functional equations implied directly by the

reflection-principle [16], the following asymptotic formula has been derived [30]

(6.3)
$$\forall k \in \mathbb{N}, \quad f_k(2n,0) \sim c_k n^{-((k-1)^2 + (k-1)/2)} (2(k-1))^{2n}, \quad c_k > 0.$$

Setting

$$w_0(x) = \frac{x^{2\sigma-2}}{1-x^2+x^{2\sigma}}$$
 and $v_0(x) = 1-x+w_0(x)x^2+w_0(x)x^3+w_0(x)x^4$

we can now state

Theorem 1. Let $k, \sigma \in \mathbb{N}$, where $k, \sigma \geq 3$, x be an indeterminate and ρ_k the dominant, positive real singularity of $\mathbf{F}_k(z)$. Then $\mathbf{T}_{k,\sigma}^{[4]}(x)$, the generating function of $\langle k, \sigma \rangle$ -structures, is given by

(6.4)
$$\mathbf{T}_{k,\sigma}^{[4]}(x) = \frac{1}{v_0(x)} \mathbf{F}_k \left(\frac{\sqrt{w_0(x)}x}{v_0(x)} \right).$$

Furthermore, the asymptotic formula

(6.5)
$$\mathbf{T}_{k,\sigma}^{[4]}(n) \sim c_k n^{-(k-1)^2 - (k-1)/2} \left(\frac{1}{\gamma_{k,\sigma}^{[4]}}\right)^n, \quad \text{for } k = 3, 4, \dots, 9.$$

holds, where $\gamma_{k,\sigma}^{[4]}$ is the minimal positive real solution of the equation $\frac{\sqrt{w_0(x)}x}{v_0(x)} = \rho_k$.

- 6.3. Combinatorics of motifs. A Motzkin-path is a path in \mathbb{Z}^2 generated by up-, down- and horizontal-steps. It starts at the origin, stays in the upper halfplane and ends on the x-axis. Let $\operatorname{Mo}_k^{\sigma}(n)$ denote the following set of Motzkin-paths:
- (a) the paths have height $\leq \sigma(k-1)$,

- (b) all up- and down-steps come only in sequences of length σ ,
- (c) all plateaux at height σ have length ≥ 3 .

Let $\mu_{k-1,\sigma}(n)$ denote the number of Motzkin-paths of length n that (a') have height $\leq \sigma(k-2)$, (b') up- and down-steps come only in sequences of length σ . We set for arbitrary $k, \sigma \geq 2$

$$G_{k,\sigma}^{*}(z) = \sum_{n\geq 0} \mu_{k,\sigma}^{*}(n)z^{n}$$

$$G_{k-1,\sigma}(z) = \sum_{n\geq 0} \mu_{k-1,\sigma}(n)z^{n}$$

$$G_{1,\sigma}(z) = \frac{1}{1-z}.$$

Now we are in position to give the main result of this section, see also Fig.21.

Proposition 3. Suppose $k, \sigma \geq 2$, then the following assertions hold:

(a) There exists a bijection

$$\beta: \mathbb{M}_k^{\sigma}(n) \longrightarrow \mathrm{Mo}_k^{\sigma}(n).$$

(b) We have the following recurrence equations

$$(6.7) \mu_{k,\sigma}^*(n) = \mu_{k,\sigma}^*(n-1) + \sum_{s=0}^{n-(2\sigma+3)} \mu_{k-1}(n-2\sigma-s)\mu_{k,\sigma}^*(s) \quad \forall n > 2\sigma+3$$

$$(6.8) \mu_{k,\sigma}(n) = \mu_{k,\sigma}(n-1) + \sum_{s=0}^{n-2\sigma} \mu_{k-1}(n-2\sigma-s)\mu_{k,\sigma}(s) \quad \forall n > 2\sigma - 1,$$

where $\mu_{k,\sigma}^*(n) = 1$ for $0 \le n \le 2\sigma + 3$ and $\mu_{k-1,\sigma}(n) = 1$ for $0 \le n \le 2\sigma - 1$.

(c) We have the following formula for the generating functions

(6.9)
$$G_{k,\sigma}^*(z) = \frac{1}{1 - z - z^{2\sigma}(G_{k-1,\sigma}(z) - (z^2 + z + 1))}$$

(6.10)
$$G_{k-1,\sigma}(z) = \frac{1}{1 - z - z^{2\sigma} G_{k-2,\sigma}(z)}.$$

In particular, for k = 3 we have the following asymptotic formula

(6.11)
$$\mu_{3,\sigma}^*(n) \sim c_{\sigma} \left(\zeta_{\sigma} \right)^{-n},$$

where c_{σ} and ζ_{σ}^{-1} are given in Tab.4.

Proof. Let \mathfrak{m} be a $\langle k, \sigma \rangle$ -motif. We construct the bijection β as follows: reading the vertex labels of \mathfrak{m} in increasing order we map each σ -tuple of origins and termini into a σ -tuple of up-steps and down-steps, respectively. Furthermore isolated points are mapped into horizontal-steps. The resulting paths are by construction Motzkin-paths of height $\leq \sigma(k-1)$. Since motifs have arcs of length ≥ 4 the paths have at height σ plateaux of length ≥ 3 . In addition we have σ -tuples of up- and down-steps. Therefore β is well defined. To see that β is bijective we construct its inverse explicitly. Consider an element $\zeta \in \mathrm{Mo}_k^{\sigma}(n)$. We shall pair σ -tuples of up-steps and down-steps as follows: starting from left to right we pair the first up-step with the first down-step tuple and proceed inductively, see Fig.21. It is clear from the definition of Motzkin-paths that this pairing procedure is well defined. Each such

pair

$$((u_i, u_{i+1}, \dots, u_{i+\sigma}, (d_j, d_{j+1}, \dots, j_{j+\sigma}))$$

corresponds uniquely to the sequence of arcs $((i+\sigma,j),\ldots,(i,j+\sigma))$ from which we can conclude that ζ induces a unique σ -canonical diagram, δ_{ζ} over [n]. Furthermore δ_{ζ} has by construction a nonnesting core. A diagram contains a k-crossing if and only if it contains a sequence of arcs $(i_1,j_1),\ldots,(i_k,j_k)$ such that $i_1 < i_2 < \cdots < i_k < j_1 < j_2 < \cdots < j_k$. Therefore δ_{ζ} is k-noncrossing if and only if its underlying path ζ has height k0. We immediately derive k1. We immediately derive k2. Whence k3 is a bijection. Using the Motzkin-path interpretation we immediately observe that k3. We constructed recursively from paths that start with a horizontal-step or an up-step, respectively. The recursions eq. (6.7) and eq. (6.8) and the generating functions of eq. (6.9) and eq. (6.10) are straightforwardly derived. As for the particular case k3. k3.

(6.12)
$$G_{3,\sigma}^*(z) = \frac{1}{1 - z - z^{2\sigma} \left[\frac{1}{1 - z - z^{2\sigma} \left[\frac{1}{1 - z} \right]} - (z^2 + z + 1) \right]}.$$

The unique dominant, real singularity of $G_{3,\sigma}^*(z)$ is a simple pole, denoted by ζ_{σ} . Being a rational function, $G_{k,\sigma}^*(z)$ admits a partial fraction expansion

$$G_{k,\sigma}^*(z) = H(z) + \sum_{(\zeta,r)} \frac{c_{(\zeta,r)}}{(\zeta - z)^r}$$

and eq. (6.11) follows in view of

(6.13)
$$[z^n] \frac{1}{\zeta - z} = \frac{1}{\zeta} [z^n] \frac{1}{1 - z/\zeta} = \frac{1}{\zeta} {n \choose 0} \left(\frac{1}{\zeta}\right)^n = \left(\frac{1}{\zeta}\right)^{n+1}.$$

6.4. **Loops.** Suppose we are given a structure S. Let α be an S-arc and denote the set of S-arcs that cross α by $\mathscr{A}_S(\alpha)$. Clearly if α cross to β , we have

$$(6.14) \beta \in \mathscr{A}_S(\alpha) \iff \alpha \in \mathscr{A}_S(\beta).$$

For two arcs $\alpha = (i, j), \alpha' = (i', j')$, we write $\alpha' \prec \alpha$ if and only if i < i' < j' < j. An arc $\alpha \in \mathscr{A}_S(\beta)$ is called a minimal, β -crossing if there exists no $\alpha' \in \mathscr{A}_S(\beta)$ such that $\alpha' \prec \alpha$. Note that $\alpha \in \mathscr{A}_S(\beta)$ can be minimal β -crossing, while β is not minimal α -crossing. We call a pair of crossing arcs (α, β) balanced, if α is minimal, β -crossing and β is minimal α -crossing, respectively. 3-noncrossing diagrams exhibit the following four basic loop-types:

(1) a hairpin-loop, being a pair

$$((i,j),[i+1,j-1])$$

where (i, j) is an arc and [i, j] is an interval, i.e. a sequence of consecutive vertices $(i, i+1, \ldots, j-1, j)$.

(2) an *interior*-loop, being a sequence

$$((i_1, j_1), [i_1 + 1, i_2 - 1], (i_2, j_2), [j_2 + 1, j_1 - 1]),$$

where (i_2, j_2) is nested in (i_1, j_1) .

(3) a multi-loop, see Fig.22, being a sequence

$$((i_1, j_1), [i_1 + 1, \omega_1 - 1], S_{\omega_1}^{\tau_1}, [\tau_1 + 1, \omega_2 - 1], S_{\omega_2}^{\tau_2}, \dots)$$

where $S_{\omega_h}^{\tau_h}$ denotes a structure containing a pseudoknot over $[\omega_h, \tau_h]$ (i.e. nested in (i_1, j_1)) and subject to the following condition: if all $S_{\omega_h}^{\tau_h} = (\omega_h, \tau_h)$, i.e. all substructures are simply arcs, for all h, then $h \geq 2$.

We finally define pseudoknots:

- (4) a pseudoknot, see Fig.23, consists of the following data:
- (P1) a set of arcs

$$P = \{(i_1, j_1), (i_2, j_2), \dots, (i_t, j_t)\},\$$

where $i_1 = \min\{i_s\}$ and $j_t = \max\{j_s\}$, such that

- (i) the diagram induced by the arc-set P is irreducible, i.e. the dependency-graph of P (i.e. the graph having P as vertex set and in which α and α' are adjacent if and only if they cross) is connected and
- (ii) for each $(i_s, j_s) \in P$ there exists some arc β (not necessarily contained in P) such that (i_s, j_s) is minimal β -crossing.

(P2) all vertices $i_1 < r < j_t$, not contained in hairpin-, interior- or multi-loops. We call a pseudoknot balanced if its arc-set can be decomposed into pairs of balanced arcs and unbalanced, otherwise, see Fig.23.

Theorem 2. Suppose $k, \sigma \geq 2$. Any $\langle 3, \sigma \rangle$ -structure has a unique loop-decomposition.

Fig.24 illustrates the loop decomposition of a 3-noncrossing structure.

6.5. Supplement of Saturation.

Suppose we are given skeleta-tree $\mathbb{T}(S_0)$ with root S_0 . Let the order of S, $\omega(S)$, denote the number of \prec -maximal S-arcs, see Fig.25. Furthermore, let $\Sigma_{i,j}$ and $\Sigma_{i,j}^{[r]}$ be some subset of structures over [i,j] and those of order r, respectively.

Let $M_{i,j}$ denote the set of saturated skeleta over [i,j] and $OSM(i,j) \in M_{i,j}$ be a mfesaturated skeleton. Furthermore, let OS(i,j) be a mfe-structure, which is a union of disjoint $OSM(i_1,j_1),...OSM(i_r,j_r)$ and unpaired nucleotides. By $OSM^{[x]}(i,j)$ and $OS^{[x]}(i,j)$ we denote the respective OSM and OS structures of order x.

In order to describe the context-sensitive saturation procedure in **cross** we denote by $OS_{mul}(i,j)$, $OS_{pk}(i,j)$ and $OS_0(i,j)$, the mfe-structures nested in a multi-loop, pseudoknot and otherwise, respectively. For a given skeleton $S_{i,j}$, we specify the mapping $S_{i,j} \mapsto OSM(S_{i,j})$ as follows: suppose $S_{i,j}$ has n_1 intervals, I_1, \ldots, I_{n_1} labelled from left to right. For given interval $I_r = [i_r, j_r]$ and $s_r \in \Sigma_{i_r, j_r}$ we consider the insertion of s_r into I_r , distinguishing four scenarios.

Case(1). I_r is contained in a hairpin-loop.

 $\underline{\omega(s_r)} = 0$. That is we have $s_r = \emptyset$. The loop generated by the s_r -insertion remains obviously a hairpin-loop, i.e. $((i_r - 1, j_r + 1), [i_r, j_r])$, with energy $H(i_r - 1, j_r + 1)$. $\underline{\omega(s_r)} = 1$. Let (p, q) be the unique, maximal s_r -arc. Then s_r -insertion produces the interior-loop

$$((i_r-1,j_r+1),[i_r,p-1],(p,q),[q+1,j_r]),$$

with energy $I(i_r - 1, j_r + 1, p, q)$. Note that $p = i_r$ implies $q \neq j_r$ and $s_r \in OSM_0^{[1]}(p,q)$.

 $\underline{\omega(s_r) \geq 2}$. In this case inserting s_r into I_r creates a multi-loop in which s_r is nested. Then $s_r \in \mathrm{OS}^{[\geq 2]}_{\mathrm{mul}}$, see Fig.26. Let $\epsilon(s)$ denote the energy of structure s. We select the set of all structures s_r such that

$$\epsilon(s_r) = \min \begin{cases} H(i_r - 1, j_r + 1) \\ I(i_r - 1, j_r + 1, p, q) + \epsilon(\mathrm{OSM}_0^{[1]}(p, q)) \\ \forall i_r \le p < q \le j_r \text{ and } p = i_r, \Rightarrow q \ne j_r \\ M + P_1 + \epsilon(\mathrm{OS}_{\mathrm{mul}}^{[\ge 2]}(i_r, j_r)). \end{cases}$$

Here, M is the energy penalty for forming a multi-loop and P_1 is the energy score of a closing-pair in multi-loop. This concludes the discussion of Case 1.

Case(2). I_r is contained in a pseudoknot.

 $\underline{\omega(s_r)=0}$. That is we have $s_r=\{\varnothing\}$ and the unpaired bases in I_r are considered to be contained in a pseudoknot.

 $\underline{\omega(s_r) \geq 1}$. In this case, s_r is a substructure which is nested in a pseudoknot, see Fig.27. As a result our selection criterion is given by

$$\epsilon(s_r) = \min \begin{cases} (j_r - i_r + 1) \cdot Q_{pk} \\ \epsilon(OS_{pk}(i_r, j_r)). \end{cases}$$

where $(j_r - i_r + 1) \in \mathbb{N}$ is the number of unpaired bases in I_r , and Q_{pk} is the energy score of the unpaired bases in a pseudoknot.

Case(3). I_r is contained in a multi-loop. In analogy to case (2), we distinguish the following cases:

 $\underline{\omega(s_r)=0}$. That is we have $s_r=\{\varnothing\}$. The unpaired bases in I_r are considered to be contained in a multi-loop.

 $\underline{\omega(s_r) \geq 1}$. In this case, s_r is a substructure nested in a multi-loop, see Fig.28. Accordingly, we select all structures satisfying

$$\epsilon(s_r) = \min \begin{cases} (j_r - i_r + 1) \cdot Q_{\text{mul}} \\ \epsilon(OS_{\text{mul}}(i_r, j_r)), \end{cases}$$

where $Q_{\rm mul}$ denotes the energy score of the unpaired bases in a multi-loop.

Case(4) I_r is contained in an interior-loop. By construction, the latter is formed by the pair (I_r, I_l) , where r < l. We then select pairs s_r in Σ_{i_r, j_r} and s_l in Σ_{i_l, j_l} . Note that only the first coordinate of the pair (I_r, I_l) is considered.

 $\underline{\omega(s_r) = 0}$ and $\underline{\omega(s_l) = 0}$. Obviously, in this case the loop formed by I_r and I_l remains an interior-loop

$$((i_r-1,j_l+1),[i_r,j_r],(j_r+1,i_l-1),[i_l,j_l]),$$

whose energy is given by $I(i_r - 1, j_l + 1, j_r + 1, i_l - 1)$.

 $\underline{\omega(s_r) \geq 1}$ and $\underline{\omega(s_l) = 0}$. In this case, $s_l = \{\emptyset\}$. I_r and I_l create a multi-loop, in which s_r and the substructure G_{j_r+1,i_l-1} are nested.

 $\underline{\omega(s_r)} = 0$ and $\underline{\omega(s_l)} \geq 1$. Completely analogous to the previous case.

 $\underline{\omega(s_r) \geq 1}$ and $\underline{\omega(s_l) \geq 1}$. In this case, I_r and I_l create a multi-loop, in which s_r , s_l and G_{j_r+1,i_l-1} are nested, see Fig.29.

Accordingly, we select all pairs of structures (s_r, s_l) satisfying

$$\epsilon(s_r) + \epsilon(s_l) = \min \begin{cases} I(i_r - 1, j_l + 1, j_r + 1, i_l - 1) \\ M + 2P_1 + \epsilon(\text{OS}_{\text{mul}}(i_r, j_r)) + (j_l - i_l + 1) \cdot Q_{\text{mul}} \\ M + 2P_1 + \epsilon(\text{OS}_{\text{mul}}(i_l, j_l)) + (j_r - i_r + 1) \cdot Q_{\text{mul}} \\ M + 2P_1 + \epsilon(\text{OS}_{\text{mul}}(i_r, j_r)) + \epsilon(\text{OS}_{\text{mul}}(i_l, j_l)) \end{cases}$$

Accordingly, we inductively saturate all intervals and in case of interior-loops intervalpairs and thereby derive $OSM(S_{i,j})$. Then we select an energy-minimal OSM(i,j)substructure from the set of all $OSM(S_{i,j})$ for any skeleton $S_{i,j}$.

7. Proofs

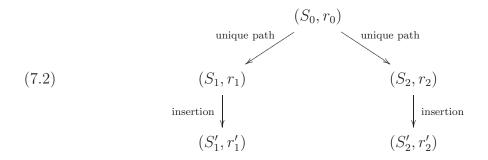
Proof of Proposition 1

Proof. We prove assertion (a) by induction on the number of inserted stacks, ℓ . For $\ell = 0$ there is nothing to prove. For $\ell = 1$, the pairs (S_0, r_0) and (S', r') differ by exactly one stack, (i, j, σ) , whence the assertion. Our objective is now to show that for any two (S'_1, r'_1) and (S'_2, r'_2) obtained from the root (S_0, r_0) via ℓ insertions, $S'_1 \neq S'_2$ holds. Suppose there exists some (\tilde{S}, \tilde{r}) , such that

(7.1)
$$(\tilde{S}, \tilde{r})$$
 insertion
$$(S'_1, r'_1) \qquad (S'_2, r'_2)$$

If the inserted stacks coincide, we have $(S'_1, r'_1) = (S'_2, r'_2)$ and there is nothing to prove. Otherwise, we obtain $S'_1 \neq S'_2$, which implies $(S'_1, r'_1) \neq (S'_2, r'_2)$, whence (a).

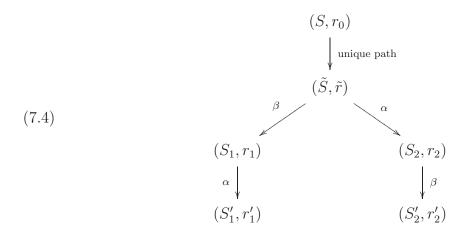
Suppose next, we have the following situation



where the uniqueness of the paths ending at (S_1, r_1) and (S_2, r_2) is guaranteed by the induction hypothesis. By assumption we have $(S_1, r_1) \neq (S_2, r_2)$ and S_1 and S_1' as well as S_2 and S_2' differ by exactly one stack. Again by induction hypothesis, we have $S_1 \neq S_2$, whence

$$(7.3) \quad (S_1, r_1) \Rightarrow_{\alpha = (i_{\alpha}, j_{\alpha}, \sigma_{\alpha})} (S'_1, r'_1), \ (S_2, r_2) \Rightarrow_{\beta = (i_{\beta}, j_{\beta}, \sigma_{\beta})} (S'_2, r'_2) \quad \text{and} \quad S_1 \neq S_2.$$

We now prove the inductive step by contradition. Suppose we have $S'_1 = S'_2$, then we can conclude that $\alpha \neq \beta$ and there exists some (\tilde{S}, \tilde{r}) such that



Indeed, we define \tilde{S} to be the skeleton derived from (S_0, r_0) by inserting all S'_1 arcs except of α, β . It is clear that the skeleton \tilde{S} exists since its stack-set is a
subset of the stack-set of S'_1 . By construction, \tilde{S} differs from S_1 and S_2 via the
stacks α and β , respectively. By induction hypothesis, there exists a unique path
from (S_0, r_0) to (\tilde{S}, \tilde{r}) , which implies the existence of a unique \tilde{r} . Furthermore, by
induction hypothesis, the paths from (S_0, r_0) to (S_1, r_1) and (S_2, r_2) are unique and
consequently contain (\tilde{S}, \tilde{r}) , whence we have the situation given in eq. (7.4).

As α and β are both minimal, without loss of generality we may assume $i_{\alpha} < i_{\beta}$. Let us consider the insertion-path $(\tilde{S}, \tilde{r}) \Rightarrow_{\beta} (S_1, r_1) \Rightarrow_{\alpha} (S'_1, r'_1)$. According to this insertion, we obtain $r_1 < i_{\alpha}$ and by construction $[r_1 + 1, i_{\beta} - 1]$ is an S_1 -interval. If $j_{\alpha} < i_{\beta}$, then α does not cross any arcs in S'_1 , which is impossible. If $j_{\alpha} > j_{\beta}$, we arrive at $\beta \prec \alpha$, which contradicts minimality of α . Therefore, we have $i_{\beta} < j_{\alpha} < j_{\beta}$, i.e. the arcs α and β are crossing. Next we consider $(\tilde{S}, \tilde{r}) \Rightarrow_{\alpha} (S_2, r_2) \Rightarrow_{\beta} (S'_2, r'_2)$. Accordingly, α must be crossed by some (\tilde{S}, \tilde{r}) -stack, say $\gamma = (i_{\gamma}, j_{\gamma}, \sigma_{\gamma})$. We next put γ into the context of the insertion-path $(\tilde{S}, \tilde{r}) \Rightarrow_{\beta} (S_1, r_1) \Rightarrow_{\alpha} (S'_1, r'_1)$ and observe that γ necessarily crosses β . Indeed, otherwise we have the following three scenarios: $i_{\gamma} > j_{\beta}$, $j_{\gamma} \leq r_1$ or $i_{\gamma} \leq r_1$, $j_{\gamma} > j_{\beta}$. In all three cases γ cannot cross α since $i_{\gamma}, j_{\gamma} \notin [r_1 + 1, i_{\beta} - 1]$, see Fig.31. As a result, γ necessarily crosses both stacks: α and β , which is a contradition to the fact that S'_1 is a 3-noncrossing skeleton, whence $S'_1 \neq S'_2$. In particular we obtain $(S'_1, r'_1) \neq (S'_2, r'_2)$, the insertion path is unique and $G_{(S_0, r_0)}$ is a tree.

In order to prove (b) we provide via Fig.30 an example, where the implication $(S_1, r_1) \neq (S_2, r_2) \Rightarrow S_1 \neq S_2$ does not hold. Note that $\mathbb{T}_{(S_0, r_0)}$ is still a tree.

Proof of Proposition 2

Proof. Let \mathscr{A}_S denote the set of S-arcs. Obviously, for any vertex $(S, r) \in \mathbb{T}(S_0)$, S is a 3-noncrossing skeleton such that $\mathscr{A}_{S_0} \subseteq \mathscr{A}_S$, whence $\mathbb{T}(S_0) \subseteq \mathbb{S}(S_0)$ holds. For an arbitrary 3-noncrossing skeleton S, let $\mathscr{A}_S^{\mathrm{ne}}$ denote the set of all nested stacks in S. Since each arc is either maximal or nested we have $\mathscr{A}_S = \mathscr{A}_{S_0} \dot{\cup} \mathscr{A}_S^{\mathrm{ne}}$. Sorting $\mathscr{A}_S^{\mathrm{ne}}$ via the linear ordering of their leftmost paired base, we obtain the sequence $\Sigma = (\alpha_1, \alpha_2, \dots, \alpha_n)$. We choose the first element $\alpha_k \in \Sigma$ which crosses some stack

in S_0 (not necessarily α_1). Then we have

$$(7.5) (S_0, i) \hookrightarrow_{\alpha_k} (S_1, r_1)$$

where, $S_1 \in \mathbb{T}(S_0)$. We proceed inductively, setting $\mathscr{A}_S^{\text{ne}} = \mathscr{A}_S^{\text{ne}} \setminus \alpha_k$ and proceed inductively until $\mathscr{A}_S^{\text{ne}} = \varnothing$. By construction, each S_k is in $\mathbb{T}(S_0)$, and $S_n = S$. Accordingly, we constructed an insertion-path in $\mathbb{T}(S_0)$ from S_0 to S, from which $\mathbb{S}(S_0) \subseteq \mathbb{T}(S_0)$ follows.

Proof of Theorem 2

Proof. Let c(S) be the core of S. We shall color the c(S)-arcs, $\alpha = (i, j)$, as follows: Case (1): $\mathscr{A}_{c(S)}(\alpha) \neq \varnothing$.

Since c(S) is a 3-noncrossing diagram, we have for any two $(i,j), (i',j') \in \mathscr{A}_{c(S)}(\beta)$, either $(i,j) \prec (i',j')$ or j < i'. Therefore for any $\beta \in \mathscr{A}_{c(S)}(\alpha)$ there exists an unique \prec -minimal arc $\alpha^* \in \mathscr{A}_{c(S)}(\beta)$ that is nested in α . If there exists some β for which $\alpha = \alpha^*(\beta)$ holds, i.e. α itself is minimal in $\mathscr{A}_{c(S)}(\beta)$, then we color α red. In other words, red arcs are minimal with respect to some crossing β . Otherwise, for any $\beta \in \mathscr{A}_{c(S)}(\alpha)$ there exists some $\alpha^*(\beta) \prec \alpha$. If $\alpha^*(\beta)$ is the unique \prec -maximal arc which is the substructure nested in α , then we color α green and blue, otherwise. Case (2): $\mathscr{A}_{c(S)}(\alpha) = \emptyset$, i.e. α is noncrossing in c(S).

If there exists no c(S)-arc $\alpha' \prec \alpha$, then we color α purple, if in the substructure nested in α there exists exactly one maximal c(S)-arc $\alpha' \prec \alpha$, we color α green and

blue, otherwise. It follows now by induction on the number of c(S)-arcs that this procedure generates a well defined arc-coloring. Let $i \in [n]$ be a vertex. We assign to i either the color of the minimal non-red c(S)-arc (r,s) for which r < i < s holds, or red if there exist only red c(S)-arcs, (r,s) with r < i < s and black, otherwise. By construction, this induces a vertex-arc coloring with the property of correctly identifying all hairpin- (purple arcs and vertices), interior- (green arcs and vertices), multi- (blue arcs and vertices) and pseudoknot (red arcs and vertices).

Acknowledgments. We are grateful to the anonymous referees who were of great help in deriving an improved version of the paper. Special thanks to Gang Ma for discussions and helpful suggestions. This work was supported by the 973 Project, the PCSIRT Project of the Ministry of Education, the Ministry of Science and Technology, and the National Science Foundation of China.

References

- [1] The pseudoknot structure data base, http://www.ekevanbatenburg.nl/PKBASE/PKB.HTML.
- [2] The HDV structure in nature, http://www.ekevanbatenburg.nl/PKBASE/PKB00075.HTML.
- [3] Mapping RNA form and function, Science 2 (2005).
- [4] T. Akutsu, Dynamic programming algorithms for RNA secondary prediction with pseudoknots, Discr. Appl. Math. **104** (2000), 45–62.
- [5] M. Andronescu, Z. C. Zhang and A. Condon, Secondary structure prediction of interacting RNA molecules, J. Mol. Biol. 345 (2005) 9871001.

- [6] S. Cao and S. J. Chen, Predicting RNA pseudoknot folding thermodynamics, Nucl. Acids. Res. 34(9) (2006), 2634–2652.
- [7] R. Cary and G. Stormo, Graph-theoretic approach to RNA modeling using comparative data, Proc. Int. Conf. Intell. Syst. Mol. Biol. 3 (1995), 75–80.
- [8] W. Y. C. Chen, E. Y. P. Deng, R. R. X. Du, R. P. Stanley, and C. H. Yan, Crossings and nestings of matchings and partitions, Trans. Am. Math. Soc. 359 (2007), 1555–1575.
- [9] W. Y. C. Chen, J. Qin, and C. M. Reidys, Crossing and nesting in tangled-diagrams, Elec. J. Comb. 15 (2008).
- [10] C. DeLisi and D. M. Crothers, Prediction of RNA secondary structure, Proc. Natl. Acad. Sci, USA 68 (1971), 2682–2685.
- [11] R. M. Dirks and N. A. Pierce, An algorithm for computing nucleic acid base-pairing probabilities including pseudoknots, J. Comput. Chem. 25 (2004), 1295–1304.
- [12] S. R. Eddy, How do RNA folding algorithms work?, Nature Biotechnology 22 (2004), 1457– 1458.
- [13] J. Edmonds, Maximum matching and polyhedron with 0,1-vertices, J. Res. Nat. Bur. Stand. 69B (1965), 125–130.
- [14] J. R. Fresco, B. M. Alberts, and P. Doty, Some molecular details of the secondary structure of ribonucleic acid, Nature 188 (1960), 98–101.
- [15] H. N. Gabow, An efficient implementation of Edmonds' algorithm for maximum matching on graphs, J. Asc. Com. Mach. 23 (1976), 221–234.
- [16] I. Gessel and D. Zeilberger, Random walk in a weyl chamber, Proc. Amer. Math. Soc. 115 (1992), 27–31.
- [17] W. Gruener, R. Giegerich, D. Strothmann, C. M. Reidys, J. Weber, I. L. Hofacker, P. F. Stadler, and P. Schuster, Analysis of RNA sequence structure maps by exhaustive enumeration i. neutral networks, Monatsh. Chem. 127 (1996), 375–389.

- [18] W. Gruener, R. Giegerich, D. Strothmann, C. M. Reidys, J. Weber, I. L. Hofacker, P. F. Stadler, and P. Schuster, Analysis of RNA sequence structure maps by exhaustive enumeration. ii, Monatsh. Chem. 127 (1996), 355–374.
- [19] C. Haslinger and P. F. Stadler, RNA Structures with Pseudo-Knots, Bull. Math. Biol. 61 (1999), 437-467.
- [20] I. L. Hofacker, Vienna RNA secondary structure server, Nucl. Acids. Res. 31(13) (2003), 3429–3431.
- [21] I. L. Hofacker, M. Fekete, C. Flamm, M. A. Huynen, S. Rauscher, P. E. Stolorz, and P. F. Stadler, Automatic detection of conserved RNA structure elements in complete RNA virus genomes, Nucl. Acids. Res. 26 (1998), 3825–2836.
- [22] I. L. Hofacker, W. Fontana, P. F. Stadler, L. S. Bonhoeffer, M. Tacker, and P. Schuster, Fast folding and comparison of RNA secondary structures, Monatsh. Chem. 125 (1994), 167–188.
- [23] J. A. Howell, T. F. Smith, and M. S. Waterman, Computation of generating functions for biological molecules, J. Appl. Math. 39 (1980), 119–133.
- [24] F. W. D. Huang, L. Y. M. Li, and C. M. Reidys, Sequence-structure relations of pseudoknot RNA, BMC Bioinformatics, 10 (2009), Suppl 1, S39.
- [25] F. W. D. Huang and C. M. Reidys, Statistics of canonical RNA pseudoknot structures, J. Theor. Biol. 253 (2008), 570–578.
- [26] M. Huynen, P. F. Stadler, and W. Fontana, Smoothness within ruggedness: the role of neutrality in adaptation, Proc. Natl. Acad. Sci, USA 93 (1996), 397–401.
- [27] E. Y. Jin, J. Qin, and C. M. Reidys, Combinatorics of RNA structures with pseudoknots, Bull. Math. Biol. 70(1) (2008), 45–67.
- [28] E. Y. Jin and C. M. Reidys, RNA-lego: Combinatorial design of pseudoknot RNA, Adv. Appl. Math. 42 (2008), 135–151.
- [29] E. Y. Jin and C. M. Reidys, Central and local limit theorems for RNA structures, J. Theor. Biol. 250(3) (2008), 547–559.

- [30] E. Y. Jin, C. M. Reidys, and R. R. Wang, Asymptoic enumeration of k-noncrossing matchings, Submitted.
- [31] I. T. Jun, O. C. Uhlenbeck, and M. D. Levine, Estimation of secondary structure in ribonucleic acids, Nature 230 (1971), 362 – 367.
- [32] D. A. M. Konings and R. R. Gutell, A comparison of thermodynamic foldings with comparatively derived structures of 16s and 16s-like rRNAs, RNA 1 (1995), 559–574.
- [33] A. Loria and T. Pan, Domain structure of the ribozyme from eubacterial ribonuclease, RNA 2 (1996), 551–563.
- [34] R. B. Lyngsø and C. N. S. Pedersen, RNA pseudoknot prediction in energy-based models, J. Comput. Biol. 7 (2000), 409–427.
- [35] G. Ma and C. M. Reidys, Canonical RNA pseudoknot structures, J. Comput. Biol. 15 (2008), 1257–1273..
- [36] D. Metzler and M. E. Nebel, Predicting RNA secondary structures with pseudoknots by mcmc sampling, J. Math. Biol. 56(1-2) (2008), 161–181.
- [37] R. Nussinov and A. B. Jacobson, Fast algorithm for predicting the secondary structure of single-stranded RNA, Proc. Natl. Acad. Sci, USA 77 (1980), 6309–6313.
- [38] J. Qin and C. M. Reidys, A combinatorial framework for RNA tertiary interaction, (2007), Submitted.
- [39] J. Reeder and G. Giegerich, Design, implementation and evaluation of a practical pseudoknot folding algorithm based on thermodynamics, Bioinformatics 5 (2004), no. 104.
- [40] C. M. Reidys, Local connectivity of neutral networks, Bull. Math. Biol. 71 (2008), 265–290.
- [41] C. M. Reidys and P. F. Stadler, Combinatorial landscapes, SIAM Review 44 (2002), 3–54.
- [42] J. Ren, B. Rastegari, A. Condon, and H. Hoos, Hotkonts: Heuristic prediction of RNA secondary structures including pseudoknots, RNA 11 (2005), 1494–1504.
- [43] E. Rivas and S. R. Eddy, A dynamic programming algorithm for RNA structure prediction including pseudoknots, J. Mol. Biol. 285(5) (1999), 2053–2068.

- [44] E. Rivas and S. R. Eddy, The language of RNA: A formal grammar that includes pseudoknots, Bioinformatics 16 (2000), 326–333.
- [45] J. Ruan, G. Stormo, and W. Zhang, An iterated loop matching approach to the prediction, Bioinformatics 20 (2004), 58–66.
- [46] P. Schuster and W. Fontana, Chance and necessity in evolution: Lessons from RNA, Physica. D. 133 (1999), 427–452.
- [47] P. Schuster, W. Fontana, P. F. Stadler, and I. L. Hofacker, From sequences to shapes and back: A case study in RNA secondary structures, Proc. Roy. Soc. Lond. B 255 (1994), 279–284.
- [48] D. B. Searls, The language of genes, Nature **420** (2002), 211217.
- [49] T. F. Smith and M. S. Waterman, RNA secondary structure, Math. Biol. 42 (1978), 31–49.
- [50] J. Tabaska, R. Cary, H. Gabow, and G. Stormo, An RNA folding method capable of identifying pseudoknots and base triples, Bioinformatics 14 (1998), 691–699.
- [51] M. Tacker, P. F. Stadler, E. G. Bornberg-Bauer, I. L. Hofacker, and P. Schuster, Algorithm independent properties of RNA secondary structure predictions, Europ. Biophy. J. 25 (1996), 115–130.
- [52] I. Tinoco, P. N. Borer, B. Dengler, M. D. Levine, O. C. Uhlenbeck, D. M. Crothers, and J. Gralla, Improved estimation of secondary structure in ribonucleic acids, Nature New Biology 246 (1973), 40–41.
- [53] Y. Uemura, A. Hasegawa, S. Kobayashi, and T. Yokomori, Tree adjoining grammars for RNA structure prediction, Theor. Comput. Sci. 210 (1999), 277–303.
- [54] M. S. Waterman, Combinatorics of RNA hairpins and cloverleaves, Stud. Appl. Math. 60 (1979), 91–96.
- [55] M. S. Waterman and T. F. Smith, Rapid dynamic programming methods for RNA secondary structure, Adv. Appl. Math. 7 (1986), 455–464.
- [56] E. Westhof and L. Jaeger, RNA pseudoknots, Curr. Opin. Struct. Biol. 2 (1992), 327–333.

[57] M. Zuker and P. Stiegler, Optimal computer folding of large RNA sequences using thermodynamics and auxiliary information, Nucl. Acids. Res. 9 (1981), 133148.

k	2	3	4	5
growth rate	2.6180	4.7913	6.8541	8.8875
k	6	7	8	9
growth rate	10.9083	12.9226	14.9330	16.9410

TABLE 1. The exponential growth rates of k-noncrossing RNA structures (minimum arc-length greater or equal than two).

k	3	4	5	6	7	8	9
$\sigma = 3$	2.0348	2.2644	2.4432	2.5932	2.7243	2.8414	2.9480
$\sigma = 4$	1.7898	1.9370	2.0488	2.1407	2.2198	2.2896	2.3523
$\sigma = 5$	1.6465	1.7532	1.8330	1.8979	1.9532	2.0016	2.0449
$\sigma = 6$	1.5515	1.6345	1.6960	1.7457	1.7877	1.8243	1.8569
$\sigma = 7$	1.4834	1.5510	1.6008	1.6408	1.6745	1.7038	1.7297
$\sigma = 8$	1.4319	1.4888	1.5305	1.5639	1.5919	1.6162	1.6376
$\sigma = 9$	1.3915	1.4405	1.4763	1.5049	1.5288	1.5494	1.5677

Table 2. Exponential growth rates of $\langle k, \sigma \rangle$ -structures.

Seq	length	TP	FP	С	CP	SP
$Bs_g lm S_P 1.1$	73	21	19	13	0.619	0.684
$HDV - It_g$	88	32	29	29	0.906	1.000
HCV_IRES	56	22	19	18	0.818	0.947
IBV_PK_NS	49	12	16	12	1.000	0.750
APLV	40	9	12	9	1.000	0.750
BCV	56	18	18	18	1.000	1.000
$BMV3_{U}PD - PK1$	26	10	10	10	1.000	1.000
PMWaV - 2	62	19	20	19	1.000	0.950

TABLE 3. Comparison of predicted and natural structures [1]. TP denotes the number of pairs in the natural structure and FP those found in cross. We also give the number of correct base pairs C, and the correct percentages of true pairs CP=C/TP and correct percentages of found pairs SP=C/FP, respectively.

σ	2	3	4	5	6	7
ζ_{σ}^{-1}	1.7424	1.5457	1.4397	1.3721	1.3247	1.2894
c_{σ}	0.1077	0.0948	0.0879	0.0840	0.0804	0.0780

Table 4. The exponential growth rates of $\mu_{3,\sigma}^*(n)$

Matrices	(i,j)	(r,s)	Matrices	(i,j)	(r,s)
whx(i,j;r,s)	unknown	unknown	vhx(i,j;r,s)	paired	paired
yhx(i,j;r,s)	unknown	paired	zhx(i,j;r,s)	paired	unknown

Table 5. The gap-matrices whx, vhx, yhx and zhx.

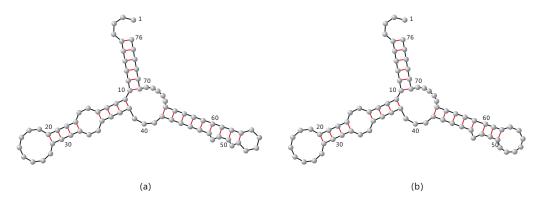


FIGURE 1. The phenylalanine tRNA: (a) represents the structure of phenylalanine tRNA, as folded by ViennaRNA [20, 22]. (b) shows the phenylalanine structure as folded by **cross** with minimum stack size 3. Note that **cross** does not contain stacks of size \leq 3, therefore (b) differs slightly from (a).

Algorithm 1

```
1: cross (seq)
 2: j \leftarrow 1
 3: while j < n do
       for i \leftarrow 0 to n - j do
 4:
          IS[i][i+j] \leftarrow Shadow (seq, i, i+j) \{Phase I\}
 5:
          Skeleton[i][i + j] \leftarrow SkeletonBranch (i, i + j) {Phase II}
 6:
          Index \leftarrow Skeleton[i][i + j]
 7:
 8:
          while Index is not empty do
            temp \leftarrow Saturation (Index) \{Phase III\}
 9:
            OSM[i][i+j] \leftarrow Optimal \text{ (temp, } OSM[i][i+j]) \text{ (Optimal: a DP-routine)}
10:
            that derives the mfe-structure.
            Index \leftarrow Index + 1
11:
          end while
12:
13:
       end for
       for i \leftarrow 0 to n - j do
14:
         for k \leftarrow i + 1 to j do
15:
            temp \leftarrow Combine (OSM[i][k]OS[k+1][i+j])
16:
            OS[i][i+j] \leftarrow Optima (OS[i][i+j], temp)
17:
18:
          end for
         OS[i][i+j] \leftarrow Optimal (OS[i][i+j], OS[i+1][i+j])
19:
       end for
20:
       j \leftarrow j + 1
21:
22: end while
23: return OS[0][n-1]
```

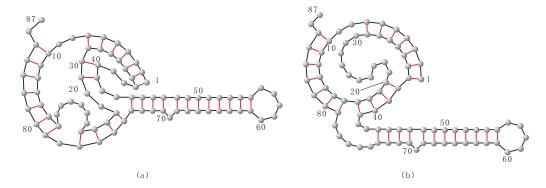


FIGURE 2. The HDV-pseudoknot structure: (a) displays the structure as folded by Rivas and Eddy's algorithm [43]. (b) shows the structure as folded by cross with minimum stack size 3.

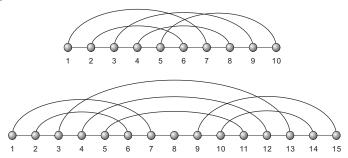


FIGURE 3. k-noncrossing diagrams: we display a 4-noncrossing, arclength $\lambda \geq 4$ and $\sigma \geq 1$ (top) and a 3-noncrossing, $\lambda \geq 4$ and $\sigma \geq 2$ diagram (bottom).

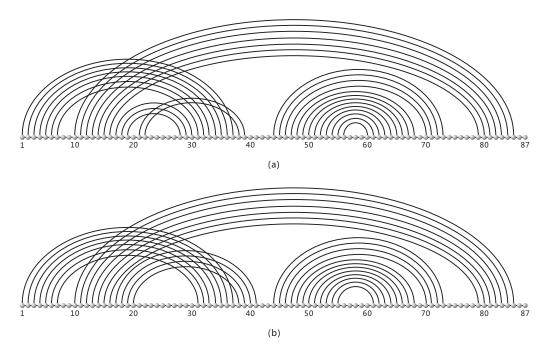


FIGURE 4. The HDV-virus pseudoknot structure: the natural structure (a) [2] versus the structure as folded by cross (b). The structure generated by cross differs from the natural structure displayed in (a) by seven base pairs.

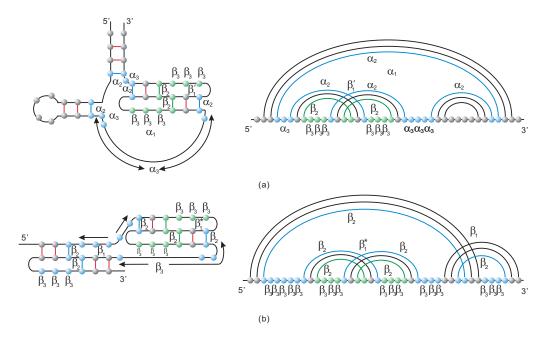


FIGURE 5. (a) a multi-loop containing a pseudoknot: as in the case of standard loops, pseudoknot base pairs contained in the multi-loop are assigned the energy contribution α_2 . The penalty for the formation of a pseudoknot within a multi-loop is given by β'_1 . (b) a pseudoknot within pseudoknot: the formation of a pseudoknot in a pseudoknot contributes β_1^{\star} .

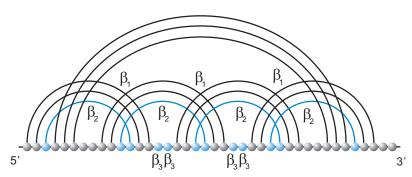


FIGURE 6. A non-planar 3-noncrossing pseudoknot and its energy $3\beta_1 + 4\beta_2 + 4\beta_3$.

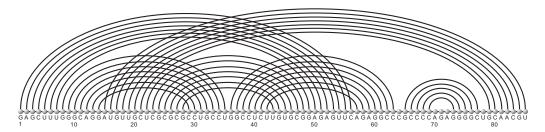


FIGURE 7. Novel pseudoknot structures: a nonplanar structure as folded by cross. As argued in Section 6.1, 3-noncrossing nonplanar structures cannot be generated iterating pairs of gap-matrices.

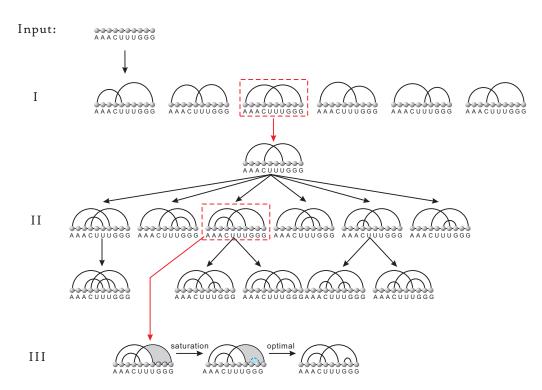


FIGURE 8. An outline of cross: the generation of shadows (Shadow), the construction of skeleta-trees (SkeletonBranch) and the filling of the skeleta (Saturation).

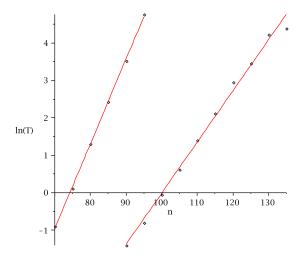


FIGURE 9. Mean folding times: we display the logarithm of the folding times of 1000 random sequences as a function of the sequence length. For 3-canonical and 4-canonical structures the linear fits are given by 0.2263n - 19.796 (left) and 0.1364n - 13.659 (right), respectively.

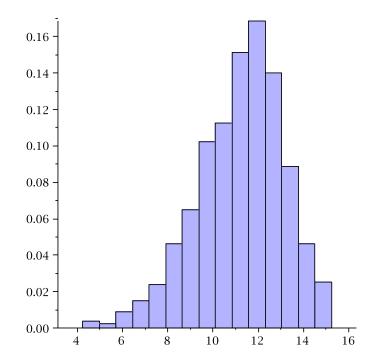


FIGURE 10. The distribution of logarithm of folding time based on 800 random sequences of length n = 100 and minimum stack length $\sigma = 3$. On the X-axis we have the logarithm of the folding time (seconds).

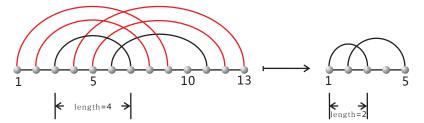


FIGURE 11. Core-structures: A structure, S, (lhs) is mapped into its core c(S) (rhs). Clearly S has arc-length ≥ 4 and as a consequence of the collapse of the stack ((4,13),(5,12),(6,11)) (the red arcs are being removed) the S-arc (3,7) is mapped into the 2-arc (1,3) in c(S).

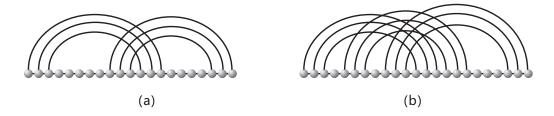


Figure 12. Motifs: a 3- and a 4-noncrossing, 3-canonical motif.

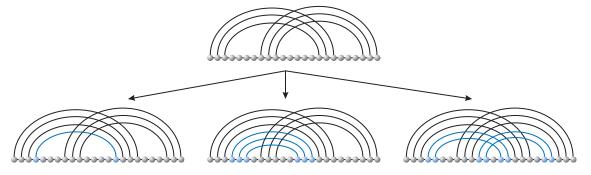


FIGURE 13. A motifs and some of its shadows: three shadows obtained from a given 3-noncrossing, 3-canonical motif.

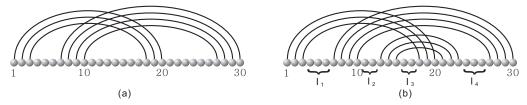


FIGURE 14. Irreducible shadows and skeleta: an irreducible shadow (a), containing the stack (1, 20, 3) and (7, 30, 4). (b) A skeleton drawn with its four induced intervals I_1, I_2, I_3, I_4 .

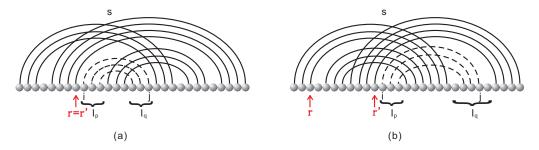


FIGURE 15. Stack-insertion: if the origin of the inserted stack (i, j, σ) is smaller than that of its predecessor (a), then r = r'. Paraphrasing the situation we can express this as "left-insertion" freezes the index r. Accordingly, (b) showcases the "right-insertion", with its induced shift of the indices $r \mapsto r'$, both indices are drawn in red.

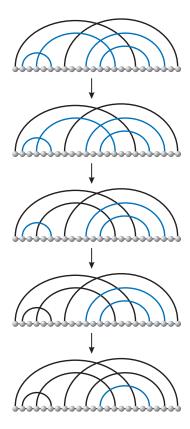


FIGURE 16. Illustration of the equality $\mathbb{T}(S_0) = \mathbb{S}(S_0)$, i.e. all skeleta are generated by the insertion procedure: we display a particular insertion path as constructed in Proposition 2.

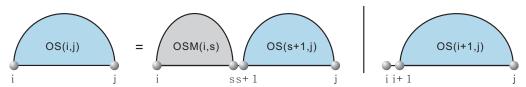


FIGURE 17. Constructing OS(i, j): inductive decomposition of the optimal structure, OS(i, j), into saturated skeleta, OSM(i, s) and unpaired nucleotides.

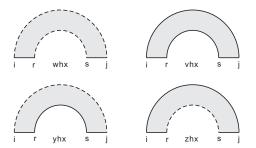


FIGURE 18. The four gap-matrices whx, vhx, yhx and zhx. The dashed line is used if the relation of two vertices is unknown, while the solid line denotes that the two vertices form a base pair.

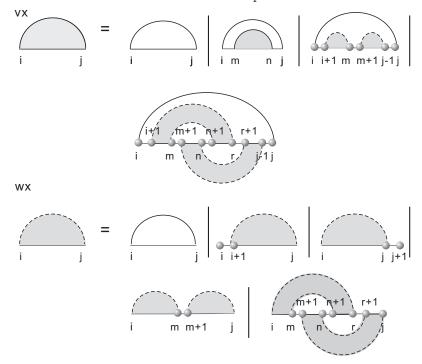


FIGURE 19. The basic recursions: recursion for vx and wx truncated at O(whx + whx + whx) in Rivas and Eddy's algorithm.

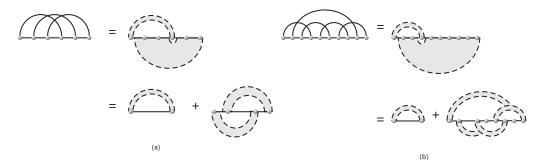


FIGURE 20. A 4-noncrossing structure which can be generated by two gap-matrices (a) and a 3-noncrossing structure, which can not be generated using two gap-matrices (b).

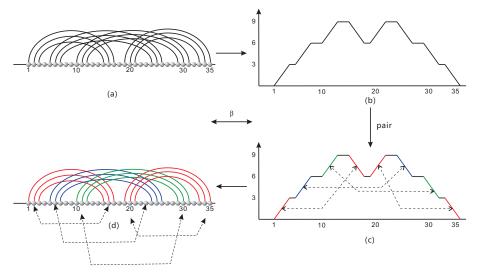


FIGURE 21. The bijection β : First we have a map from (a) to (b). Then we pair the σ -tuples of up-steps and down-steps, see the vertical map from (b) to (c). The so derived pairs, see the horizontal map from (c) to (d), allow to reconstruct the original motif.

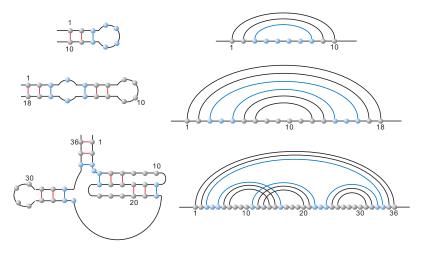


FIGURE 22. The standard loop-types: hairpin-loop (top), interior-loop (middle) and multi-loop (bottom). These represent all loop-types that occur in RNA secondary structures.

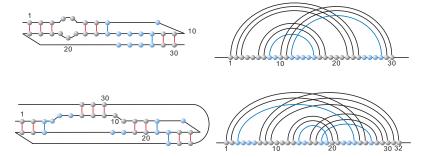


FIGURE 23. Pseudoknots: we display a balanced (top) and an unbalanced pseudoknot (bottom). The latter contains the stack over (3,24), which is minimal for the arc (9,30), which is not contained in the pseudoknot.

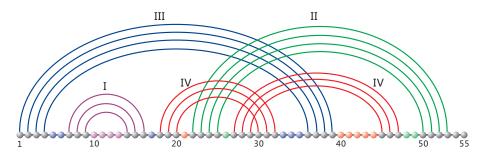


FIGURE 24. Shadows and loops: we give the loop-decomposition illustrating Theorem 2. Here I (purple) is a hairpin-loop, II (green) represents an interior-loop, III (blue) is a multi-loop and finally IV (red) is a (balanced) pseudoknot.

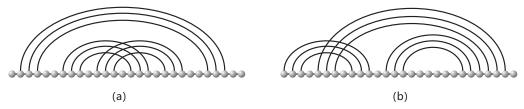


FIGURE 25. Order: In (a) we display a structure of order one. (b) show-cases a structure of order two.

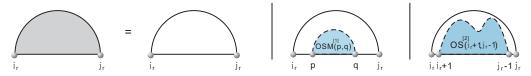


FIGURE 26. Saturation in hairpin-loops: the interval on the left hand side is filled with substructures s_r such that $\omega(s_r) = 0$ (left), $\omega(s_r) = 1$ (middle) or $\omega(s_r) \geq 2$ (right).

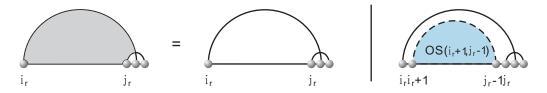


Figure 27. Saturation of an interval nested in a pseudoknot.

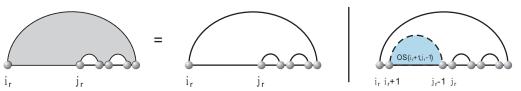


Figure 28. Saturation of an interval contained in a multi-loop.

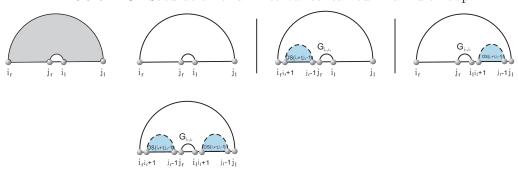


FIGURE 29. Saturation of an interval contained in an interior-loop, which is obtained by I_r and I_l , where r < l.

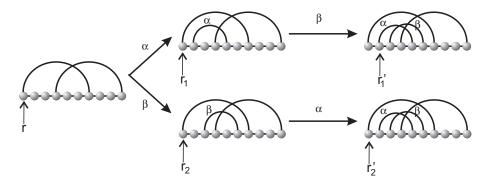


FIGURE 30. Illustration of assertion (b) of Proposition 1: the case k > 3. While $\mathbb{T}_{(S_0,r_0)}$ is still a tree (over pairs), the implication $(S_1,r_1) \neq (S_2,r_2) \Rightarrow S_1 \neq S_2$ does not hold in general.

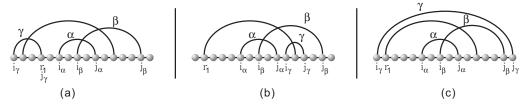


FIGURE 31. Illustration of the proof of Proposition 1. The three different scenarios for a noncrossing γ , representing stacks by isolated arcs. (a) $j_{\gamma} \leq r_1$, (b) $i_{\gamma} > j_{\beta}$ and (c) $i_{\gamma} \leq r_1$, $j_{\gamma} > j_{\beta}$.