# RNA Triplet Repeats: Improved Algorithms for Structure Prediction and Interactions

## Kimon Boehmer
Laboratoire d'Informatique de l'Ecole Polytechnique (LIX UMR 7161), Institut Polytechnique de Paris, France

## Sarah J. Berkemer 🏠 ⓘ
Laboratoire d'Informatique de l'Ecole Polytechnique (LIX UMR 7161), Institut Polytechnique de Paris, France

## Sebastian Will 🏠 ⓘ
Laboratoire d'Informatique de l'Ecole Polytechnique (LIX UMR 7161), Institut Polytechnique de Paris, France

## Yann Ponty[1] ✉ ⓘ
Laboratoire d'Informatique de l'Ecole Polytechnique (LIX UMR 7161), Institut Polytechnique de Paris, France

## ──── Abstract ────

RNAs composed of Triplet Repeats (TR) have recently attracted much attention in the field of synthetic biology. We study the mimimum free energy (MFE) secondary structures of such RNAs and give improved algorithms to compute the MFE and the partition function. Furthermore, we study the interaction of multiple RNAs and design a new algorithm for computing MFE and partition function for RNA-RNA interactions, improving the previously known factorial running time to exponential. In the case of TR, we show computational hardness but still obtain a parameterized algorithm. Finally, we propose a polynomial-time algorithm for computing interactions from a base set of RNA strands and conduct experiments on the interaction of TR based on this algorithm. For instance, we study the probability that a base pair is formed between two strands with the same triplet pattern, allowing an assessment of a notion of orthogonality between TR.

---

[1] To whom correspondence should be addressed

24th International Workshop on Algorithms in Bioinformatics (WABI 2024).
Editors: Solon P. Pissis and Wing-Kin Sung; Article No. 18; pp. 18:1–18:23
Leibniz International Proceedings in Informatics
LIPICS Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

## 1 Introduction

RNAs composed of Triplet Repeats (TR) have attracted much attention, and harbour promises in the field of synthetic biology, due to their demonstrated capacity to self-assemble into droplets [13, 10]. Those can in turn be used to compartmentalize cellular processes, thereby creating a "clean room", free of the natural cellular clutter, where synthetic circuits can be executed without interference. The exact process underlying this phenomena is still the object of ongoing investigations, but it is hypothesized that repetitive RNAs may induce Liquid-Liquid Phase separation mediated by unstable/transient structures. Repetitive RNAs are also found at the origin of severe Neurological Triplet Expansion Diseases (TED), including Friedreich attaxia [21] and Triplet Repeat Diseases (TRD) such as Huntington disease [14]. For multiple TEDs and TRDs, overly expanded RNAs have been observed to aggregate into RNA foci, leading to a sequestration of RNA binding proteins. Local secondary structures and interactions are impacted by the repeat, and generally believed to contribute to the pathogenicity and treatment efficiency. To study those phenomena *in silico*, and in particular the impact of the repeated motif and number of repeats on aggregates, one needs to predict the MFE structure of potentially large RNAs, and many-body interactions. Recently, coarse-grained simulations showed a disparity between odd or even numbers of triplet repeats [17] as well as extensions to quadruplet and non-redundant tandem repeats [1].

RNA folding by energy minimization is a classic algorithmic problem in Bioinformatics, historically solved in time $\Omega(n^3)$ using dynamic programming [19, 23]. Despite recent suggestions for heuristics [12], the best algorithm to date to solve energy minimization has runtime $\mathcal{O}(n^{2.8603})$ [5], and both its implementation and extension beyond a base-pair maximization setting represent considerable challenges. Prior works have also investigated conditional lower bounds, and found that the existence of a $\mathcal{O}(n^{2-\varepsilon})$ algorithm would refute the Strong Exponential Time Hypothesis (SETH) [5]. Meanwhile, an $\mathcal{O}(n^{\omega-\varepsilon})$ algorithm would disprove the $k$-clique conjecture, with $\omega < 2.373$ being the matrix multiplication exponent [5, 6].

RNA-RNA interaction prediction represents an equally relevant, yet computationally substantially more involved algorithmic problem. For a fixed number of interacting strands, polynomial-time algorithms have been proposed. For example, by excluding so-called zig-zag joint conformations, Alkan et al. [2] proposed a polynomial-time algorithm for the interaction of two strands, while also showing **NP**-hardness for the case where we include these conformations. In the unbounded case, Dirks et al. [9] gave a factorial-time algorithm for computing the partition function (PF) over multiple strands. Additionally, it was shown that energy minimization in this setting is **APX**-hard (and by that **NP**-hard) [7], even for a very simple energy model.

In this work, we show that the repeated nature of RNA can be exploited to obtain substantially improved algorithms for several problems. First, we show that the MFE of a triplet-repeat RNA can be predicted in linear time, both with respect to base pair maximization and Turner energy model, and is realized by either the open chain or a single helix. We then consider the interaction of multiple triplet repeats and propose improved algorithms for the general (non-triplet) case as well as algorithms specifically for the interaction of TR. For the latter case, we show **NP**-hardness in a reasonable energy model. We then propose a polynomial-time algorithm for the setting where we are given a "soup" of strands instead of a fixed set, and, using this algorithm, conduct experiments on the probability that a base pair is folding, interacting with another identical sequence or interacting with a different sequence.

## 2    Definitions and Problems Statement

### 2.1    Definitions

**RNA sequence and folding.**    An RNA sequence (or just sequence) is a word $s \in \{A, C, G, U\}^+$. The length of $s$ is denoted by $|s|$ and the $i$-th position of $s$ by $s_i$. A position on a sequence is also called a base. We associate to each base $s_i$ its letter by $l(s_i)$. We define $P :=$ $\{\{C, G\}, \{A, U\}, \{G, U\}\}$. A (pseudoknot-free) secondary structure $S$ is a set of unordered pairs of bases, hereunder called base pairs, such that:

- each base pair is a Watson-Crick or Wobble pair, i.e. for all $\{s_i, s_j\} \in S$, $\{l(s_i), l(s_j)\} \in P$;
- each base is involved in at most one base pair, i.e. for all bases $s_i$, $|\{p \in S \mid s_i \in p\}| \leq 1$;
- $S$ is pseudoknot-free, i.e. there are no $\{s_i, s_j\}, \{s_k, s_\ell\} \in S$ with $i < k < j < \ell$;
- each base pair encloses at least $\theta$ bases, i.e. if $\{s_i, s_j\} \in S$, then $j - i > \theta$. We usually call $\theta$ the minimal base pair span, and use $\theta = 3$ unless explicitly specified.

We denote by $\Omega(s)$ or $\Omega$ the set of all pseudoknot-free secondary structures over sequence $s$.

We associate each secondary structure $S \in \Omega$ to a free energy, according to an energy model $E : \{A, C, G, U\}^+ \times \Omega \to \mathbb{R}$. For example, in the base pair model $E_{\mathrm{bp}}$, we simply count the number of base pairs in $S$, hence set $E_{\mathrm{bp}}(s, S) = -|S|$. More advanced energy models reason about the free energy introduced by motifs occurring in the secondary structure, such as the loops considered by the Turner nearest-neighbor model [22].

**Interactions.**    A strand is an RNA sequence which is identified as a unique object in a set. In other words, in a set of strands $R$, we can have two strands $s \neq r$ that consist of the same sequences, that is $l(s_i) = l(r_i)$ for all $i \in \{1, ..., |s| = |r|\}$, but still are different objects. To describe the interaction of multiple strands, we are given a set $R$ of strands, where $m := |R|$.

A *circular permutation* $\pi : R \to \{0, ..., m - 1\}$ of a strand set $R$ is a permutation of all elements in $R$ except for one fixed strand $s^*$, which is fixed to position 0. Then, the bases are naturally ordered by $s_i <_\pi r_j \equiv s < r \vee (s = r \wedge i < j)$. We define $O_\pi$ as the set of all tuples of bases $(s_{i_1}^1, ..., s_{i_k}^k)$ such that there is a $j$ with $s_{i_j}^j <_\pi s_{i_{j+1}}^{j+1} <_\pi ... <_\pi s_{i_k}^k <_\pi s_{i_1}^1 <_\pi ... <_\pi s_{i_{j-1}}^{j-1}$.
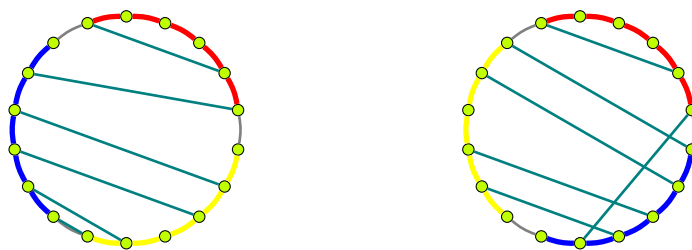
A *secondary structure* $S$ of a strand set $R$ is a set of base pairs $\{s_i, r_j\}$ from strands in $s, r \in R$ such that $\{l(s_i), l(r_j)\} \in P$, each base appears in at most one base pair and each intra-strand base pair encloses at least $\theta$ bases, i.e. $\{s_i, s_j\} \in S \to j - i > \theta$.

The *polymer graph* of a secondary structure $S$ and a circular permutation $\pi$ on $R$ is a graph $G = (V, E)$ with $V := \{s_i \mid s \in R, 1 \leq i \leq |s|\}$ and $E := S \cup \{\{s_i, s_{i+1}\} \mid s \in R, 1 \leq i < |s|\} \cup C := \{\{s_{|s|}, r_1\} \mid (\pi(s) + 1) \bmod |R| = \pi(r)\}$. The edges $E - S$ are drawn in a cycle (naturally induced by the circular permutation), while the edges in $S$ are drawn as straight lines between the bases. Examples for the polymer graphs of a single secondary structure under two different circular permutations can be found in Figure 1.

Two strands $s, r$ are connected if there is a path from $s_1$ to $r_1$ that does not use edges from $C$. A secondary structure is connected if all of its strands are connected. Note that connectedness is independent of the circular permutation $\pi$.

A secondary structure $S$ is called *pseudoknot-free* if there is a circular permutation $\pi$ such that there are no crossing lines in the polymer graph, or formally, there are no two base pairs $\{s_i, t_k\}, \{u_\ell, r_j\} \in S$ with $(s_i, u_\ell, t_k, r_j) \in O_\pi$. The set of all pseudoknot-free secondary structures over a strand set $R$ is denoted by $\Omega(R)$.

As for the folding, we associate to each $S \in \Omega(R)$ a free energy $E : 2^{\{A,C,G,U\}^*} \times \Omega \to \mathbb{R}$. In the base pair model, apart from the number of base pairs $p$ of base pairs, we also add a strand association penalty $K_{\mathrm{assoc}}$ for each of the $(m - \ell)$ strand associations, where $\ell$ is the number of connected components (also called complexes) of $S$. Thus, the free energy of $S \in \Omega$ in this model is defined as $E(R, S) = -p + (m - \ell)K_{\mathrm{assoc}}$.

■ **Figure 1** The same secondary structure on a strand set with three strands drawn in two different circular permutations. The strands are depicted by the blue, red and yellow lines while green lines indicate base pairs. Gray lines connect subsequent strands and depend on the strand permutation.

## 2.2  Computational problems

For a single strand, the two classical problems in RNA bioinformatics are:

---
MINIMUM FREE ENERGY (MFE) UNDER ENERGY MODEL $E$
**Input:** A sequence $s$
**Output:** Minimum free-energy $\min_{S \in \Omega(s)} E(s, S)$

---

---
PARTITION FUNCTION UNDER ENERGY MODEL $E$
**Input:** A sequence $s$ and a positive temperature $T$ in Kelvin (K)
**Output:** Partition function $\mathcal{Z}_s := \sum_{S \in \Omega(s)} \exp\{\frac{-E(s,S)}{kT}\}$

---

where $k = 1.987 \cdot 10^{-3}$kcal.mol$^{-1}$.K$^{-1}$ is the Boltzmann constant.

In the multi-strand setting, we focus on energy minimization. In Dirks et al. [9], the authors adopt a thermodynamic perspective on the free energy of a secondary structure over multiple strands, such that potential rotational symmetries require an adjustment of the computed value. For the MFE, we focus on a more algorithmic perspective, where all rotationally symmetric structures are elements of a search space, and a simple base pair energy model. In our main algorithmic problem of interest, we are given a set of strands and are looking for the minimum free energy of the secondary structure over these strands:

---
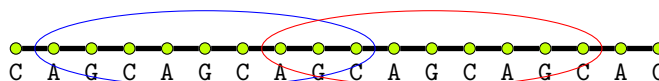MFE STRAND INTERACTION
**Input:** Set of strands $R_0$
**Output:** $\min_{S \in \Omega(R_0)} E(R_0, S)$

---

We also consider a slightly different setting, where the number of occurrences of each triplet/strand is unconstrained beyond the total number $m$ of interacting strands. This allows to study situations where the strands concentrations are in excess, so that sequences can be locally seen as infinitely available often within a set (or "soup") $R$ of strands. We then look for the best structure over $m$ strands that all appear in $R$. Each sequence in the soup can appear zero or multiple times in a secondary structure. More formally:

---
MFE STRAND SOUP INTERACTION
**Input:** Set of *sequences* $R = \{r_1, ..., r_p\}$, $m \in \mathbb{N}$ encoded in unary
**Output:** $\min_{t_1 \in R, ..., t_m \in R} \min_{S \in \Omega(\{t_1, ..., t_m\})} E(\{t_1, ..., t_m\}, S)$

---

**Figure 2** The blue and red region of the TR sequence are identical.

## 2.3 Triplet repeats RNAs and their properties

**Triplet repeat RNAs (TR).** Of special interest to us are RNA sequences that are composed of *triplet repeats* (TR), that is, they have the form $(X \cdot Y \cdot Z)^k$ for $X, Y, Z \in \{A, C, G, U\}$ and $k \in \mathbb{N}^+$. We will describe how we can improve the general algorithms for the above computational problems in the case of TR.

An algorithmically convenient property about a region $[s_i, s_j]$ of a TR sequence is:

▶ **Observation 1.** *For a triplet repeat sequence $s$ and $1 \le i \le j \le |s|$,*

$$[s_i, s_j] = [s_{i \bmod 3}, s_{j-(i-i \bmod 3)}].$$

In other words, we can shift any region three positions to the left or right, and in particular we can shift it to the beginning of the sequence, as visualized in Figure 2. That way, the index that usually denotes the beginning of the considered sequence in a dynamic programming (DP) algorithm can be restricted to values 1, 2 and 3. Hence, the length of the value range is constant and not linear anymore, which gives an easy linear improvement of running time and storage for MFE as well as PF computation.

We also note that TR sequences can be encoded exponentially more compact than general sequences. Each TR sequence is uniquely identified by its pattern $XYZ \in \{A, C, G, U\}^3$ and its number of repeats $k$. In other words, $6 + \lceil \log_2 k \rceil$ bits are enough to encode a TR sequence with $k$ repeats. We will refer to this encoding as the *compact* encoding, while the *explicit* encoding consists of the complete sequence $s \in \{A, C, G, U\}^{3k}$ (the latter can also be seen, asymptotically equivalent, as a compact encoding where $k$ is encoded in unary).

Looking into more structural properties of triplet repeats, we can observe that, since each base repeats after two other bases, there cannot be a base pair that encloses exactly 2 bases. Thus, requiring two ($\theta = 2$) or three ($\theta = 3$) enclosed bases between any base pair is equivalent:

▶ **Observation 2.** *A secondary structure $S$ for $(XYZ)^k$ fulfills minimum base pair span $\theta$ with $\theta \equiv_3 2$ if and only if it fulfills minimum base pair span $\theta + 1$.*

Finally, if we consider the graph $G = (\{A, C, G, U\}, P)$, where $P$ is the set of allowed base pairs, we can see that it does not contain any triangles. From this we can observe:

▶ **Observation 3.** *For any triplet sequence $(XYZ)^k$, there is a letter $V \in \{X, Y, Z\}$, that we call the **covering letter**, that is contained in all base pairs, i.e. $V \in p$ for all $p \in S$ and $S \in \Omega$.*

## 3 Single-Stranded Triplet Repeats

Our goal is to specify the exact MFE, and the corresponding secondary structure, when given a triplet pattern $XYZ$ and length $k$ of our TR sequence $s$, as well as the minimum base pair span $\theta$. This will give us a very efficient way of computing the MFE in this simple setting.

## 3.1    Linear time solution for base pair maximization

We first consider the properties of the MFE structure for TR RNAs in a base pair maximization model, where the free energy $E_{\mathrm{bp}}$ of a secondary structure $S \in \Omega$ is such that $E_{\mathrm{bp}}(s, S) = -|S|$.

    We can first prove an upper bound on the number of base pairs in a TR sequence:

▶ **Lemma 4.** *Consider a TR sequence* $s := (XYZ)^k$ *and a minimum number of enclosed bases* $\theta \geq 0$, *such that* $\lfloor \frac{\theta+1}{3} \rfloor \leq k$. *We have* $E_{bp}(s, S) \leq k - \lfloor \frac{\theta+1}{3} \rfloor$ *for any* $S \in \Omega(s)$.

**Proof.** Without loss of generality, let $Z$ be the covering letter of $s$. Any non-empty secondary structure has an innermost base pair which must respect the minimum base pair span $\theta$. For $\theta = 2$, which is equivalent to $\theta = 3$ by Observation 2, as well as for $\theta = 4$, at least one $Z$ base must remain unpaired, and increasing $\theta$ by 3 will result into one new unpairable $Z$ base. Thus we know that at least $\lfloor \frac{\theta+1}{3} \rfloor$ $Z$ bases will remain unpaired and at most $k - \lfloor \frac{\theta+1}{3} \rfloor$ $Z$-bases are pairable. Since every base pair must involve a $Z$ base, we can conclude.     ◀

We now show that this upper bound is almost always tight. To this end, first notice that for all triplet patterns $XYZ$ such that $\{\{X,Y\}, \{X,Z\}, \{Y,Z\}\} \cap P = \emptyset$, no base pair can be built and thus the maximum value is trivially 0. We call TR sequences of such patterns non-folding, and all other TR sequences folding.

▶ **Lemma 5.** *For* $\theta \in \{0,1\}$ *and* $k > 1$, *we always have* $E(s, S) = k$ *for any secondary structure* $S$ *over a folding sequence* $s = (XYZ)^k$.

**Proof.** If $\{X, Z\} \in P$, connect $X$ and $Z$ in each triplet. Else, connect the outermost pair (say without loss of generality $\{X, Y\}$). We obtain the inner sequence $(YZX)^{k-1}$ (with $k - 1 > 0$) and we can proceed as above since $\{Y, X\} \in P$.     ◀

For the more natural case $\theta > 1$, the upper bound from Lemma 4 is not always tight. The next lemma exactly specifies the MFE and its structure:

▶ **Lemma 6.** *Let* $\theta > 1$. *The minimum MFE structure of a folding sequence* $(XYZ)^k$ *has value*
- $k - 1 - \frac{\theta-1}{3}$, *if* $(\{X, Z\} \notin P \wedge (\theta + 3k) \equiv_6 4) \vee (\{X, Y\}, \{Y, Z\} \notin P \wedge (\theta + 3k) \equiv_6 1)$
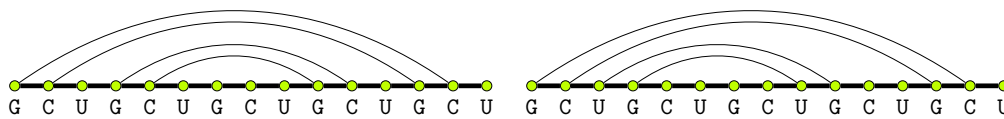- $k - \lfloor \frac{\theta+1}{3} \rfloor$, *otherwise*

*Furthermore, a minimum MFE structure is obtained by a single helix of base pairs of one letter pair* $p$. *If both* $\{X, Z\} \in P$ *and one of* $\{X, Y\}$ *and* $\{Y, Z\} \in P$, *we set* $p := \{X, Z\}$ *if* $(\theta + 3k) \equiv_6 4$ *and* $p := \{X, Y\}$ *(or* $p := \{Y, Z\}$) *if* $(\theta + 3k) \equiv_6 1$; *otherwise, we set* $p$ *to the letters of an arbitrary pairable base pair.*

The proof of this lemma involves many case distinctions and can be found in the appendix. Setting $\theta = 3$, we get the following corollary:

▶ **Corollary 7.** *In the base pair maximization model, if* $\theta = 3$, *the MFE structure of any TR sequence* $(XYZ)^k$ *has* $k - 1$ *base pairs.*

    Determining the MFE is thus a simple calculation taking logarithmic time in the (explicit) size of the triplet repeat sequence. From this we can derive:

▶ **Theorem 8.** *MFE prediction for compactly encoded TR in the base pair maximization model can be solved in linear time.*

**Figure 3** Two different optimal secondary structures for $\texttt{GCU}_5$, for $\theta = 3$.

▶ **Remark 9.** The optimal secondary structure does not need to be unique. In particular, for a simple energy model, the number of optimal secondary structures for triplet repeats can even be exponential. For example, consider the sequence $(\texttt{GCU})^k$ as illustrated in Figure 3. When constructing the base pairs from outside to inside, in every step, we can choose whether we add the base pairs $\texttt{G-U}, \texttt{U-G}$ or the base pairs $\texttt{G-C}, \texttt{C-G}$. This decision can be repeated $\lfloor \frac{k}{2} \rfloor - 1$ times (assuming $\theta = 3$), giving $\Omega(2^{k/2})$ different optimal secondary structures.

## 3.2 Minimum Free-energy in the Turner model

For the Turner model, we will argue that the optimal structures obtained for BP maximization remains optimal for the Turner nearest neighbor model under reasonable assumptions, satisfied by current versions of the model [22]. We first show a helpful lemma:

▶ **Lemma 10.** *Assume that a TR region of $s$ where the covering letter appears $k$ times has $B$ branches. the number of base pairs is at most $k - B$.*

**Proof.** Let $V$ be the covering letter of $s$. By Observation 3, for each base pair $\{s_i, s_j\}$, either $l(s_i) = V$ or $l(s_j) = V$. Furthermore, each of the $B$ branches contains one unpairable $V$-base (since $\theta = 3$). Thus, there are only $k - B$ pairable $V$-bases and we conclude. ◀

We show the absence of multiloops, *i.e.* structural motifs consisting of $B \geq 2$ branches, in the Turner MFE, with some simplifications. Their free energy contribution is composed of an initiation penalty $\alpha$, a value $\beta$ for each branch, and an asymmetry penalty $\gamma$. The overall contribution of a multiloop $S$ is given by $E(s, S) = \alpha + \beta B + \gamma C + E_{\text{in}}$, where $E_{\text{in}}$ is the MFE of the interior secondary structure of the branches. We will assume $N := \min_{V,W \in \{X,Y,Z\} : \{V,W\} \in P} E_{V,W}$ to be the best contribution of a single base pair appearing in a stacking in our triplet pattern, and we will not consider dangling ends etc.

▶ **Lemma 11.** *Any Turner-MFE secondary structure $S^*$ over a TR sequence does not contain any mutliloops, assuming $\beta \geq N, \alpha > -\beta, \gamma \geq 0$.*

**Proof.** Let $S$ be a multiloop structure on region $s$ with $k$ appearances of the covering letter and let $S^*$ be a stacking on the same region. Their free energy values are related as follows:

$$E(s, S) \geq \alpha + \beta B + \gamma C + (k - B)N \tag{1}$$
$$> -\beta + \beta B + (k - B)N \tag{2}$$
$$= (k - 1)N + (\beta - N)(B - 1) \tag{3}$$
$$\geq (k - 1)N \tag{4}$$
$$\geq E(s, S^*) \tag{5}$$

where (1) comes from our above observation and Lemma 10, (2) from $\alpha > -\beta$ and $\gamma \geq 0$, (4) from $\beta \geq N$ and $B \geq 2$ (by definition of a multiloop). For inequality (5), first notice that $S^*$ contains $k - 1$ base pairs by Corollary 7. As noticed in Remark 9, we can choose which base pair is used in $S^*$ without affecting the optimality. In particular, we can always choose the base pair consisting of the letters $V, W$ that optimize their contribution, such that $E_{V,W} = N$. We get $E(s, S^*) \leq (k - 1)N$. ◀

▶ **Remark 12.** The above assumptions are satisfied by the Turner 2004 energy model ($\alpha = 9.25$, $\beta = -0.63$, $\gamma = 0.91$ and $N \leq -0.93$) [22].

Lemma 10 also excludes secondary structures with multiple exterior faces. Thus, by the above two lemmata, we can conclude that the MFE in the Turner model is also of the canonical form described in the BP maximization setting.

## 3.3 Linear-time computation of the partition function

In the context of computing the partition function, one can write a weighted context-free grammar which, for any given pattern $XYZ$, simultaneously generates all TR sequences along with their associated set of secondary structures $\Omega$.

Below is the context-free grammar for the pattern `CAG`:

$$
\begin{aligned}
S_C^G \rightarrow{} & ( \cdot_A\, S_G^C\, \cdot_A\, ) & | \, ( \cdot_A\, S_G^C\, \cdot_A\, )\, S_C^G & & | \, \cdot_C\, \cdot_A\, S_G^G & & | \, \cdot_C\, \cdot_A\, \cdot_G \\
S_G^C \rightarrow{} & (\, S_C^G\, ) & | \, (\, S_C^G\, )\, \cdot_A\, S_G^C & & | \, \cdot_G\, S_C^G & & | \, \cdot_G\, \cdot_C \\
S_G^G \rightarrow{} & (\, S_C^G\, )\, \cdot_A\, \cdot_G & | \, (\, S_C^G\, )\, \cdot_A\, S_G^G & & | \, \cdot_G\, S_C^G & & \\
S_C^C \rightarrow{} & ( \cdot_A\, S_G^C\, \cdot_A\, )\, \cdot_A & | \, ( \cdot_A\, S_G^C\, \cdot_A\, )\, S_C^C & & | \, \cdot_C\, \cdot_A\, S_G^C & &
\end{aligned}
$$

Namely, the terminal $S_C^G$ generates all secondary structures for the RNA sequence $(CAG)^k$ for all $k > 0$, $S_G^C$ the structures of $(GCA)^k GC$ for $k \geq 0$, $S_G^G$ the structure of $G(CAG)^k$ for $k > 0$, and $S_C^C$ corresponds to the pattern $(CAG)^k C$ for some $k > 0$.

Following standard methodologies in enumerative/analytic combinatorics [8], such a grammar can be generically translated into a system of functional equations involving weighted generated functions for each non-terminal:

$$
\begin{aligned}
S_C^G(z) &= \beta\, z^4\, S_C^G(z) + \beta\, z^4\, S_G^C(z)\, S_C^G(z) + z^2\, S_G^G(z) + z^3 \\
S_G^C(z) &= \beta\, z^2\, S_C^G(z) + \beta\, z^3\, S_C^G(z)\, S_G^C(z) + z\, S_C^G(z) + z^2 \\
S_G^G(z) &= \beta\, z^4\, S_C^G(z) + \beta\, z^3\, S_C^G(z)\, S_G^G(z) + z\, S_C^G(z) \\
S_C^C(z) &= \beta\, z^3\, S_G^C(z) + \beta\, z^2\, S_G^C(z)\, S_C^C(z) + z^2\, S_G^C(z)
\end{aligned}
$$

where $\beta := e^{1/kT}$ is the Boltzmann weight associated to base pairs and, in particular:

$$
S_C^G(z) = \sum_{s \in \mathcal{L}(S_C^G)} \beta^{\#\mathrm{BP}(s)}\, z^{|s|} = \sum_{k \geq 0} \sum_{\substack{s \in \mathcal{L}(S_C^G) \\ \text{such that } |s| = 3\,k}} e^{\frac{\#\mathrm{BP}(s)}{kT}}\, z^{3k} = \sum_{k \geq 0} \mathcal{Z}_{(CAG)^k}\, z^{3k}
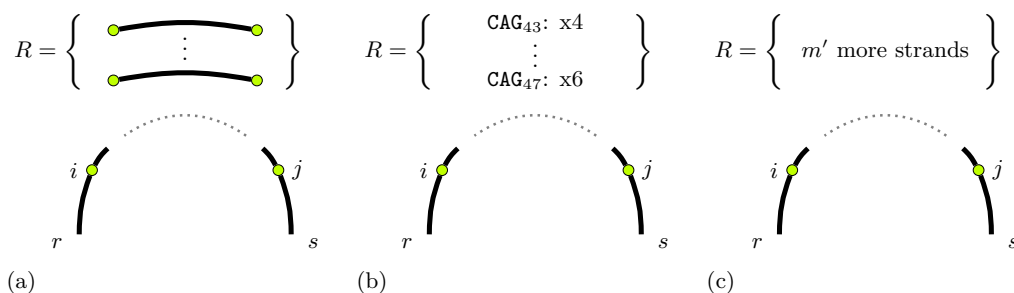$$

The partition function of $\mathcal{Z}_{(CAG)^k}$ can then be obtained as $[z^{3k}]\, S_C^G(z)$, the coefficient of degree $3k$ in $S_C^G(z)$. Since the system of functional equations is algebraic, the coefficients of each generating function obey a linear recurrence with polynomial coefficients [15], which can be efficiently [4] and effectively computed [20]. We obtain an equation of the form:

$$
\mathcal{Z}_{(CAG)^k} = P_1(k)\, \mathcal{Z}_{(CAG)^{k-1}} + P_2(k)\, \mathcal{Z}_{(CAG)^{k-2}} + \cdots + P_d(k)\, \mathcal{Z}_{(CAG)^{k-d}}
$$

where each $P_i$ is a polynomial in $k$, and $d$ is a constant . $\mathcal{Z}_{(CAG)^k}$ can then be computed using a linear number of arithmetic operations. This also holds for other triplets and thus:

▶ **Theorem 13.** *The partition function of a TR can be computed in $\Theta(k)$ arithmetic operations.*

**Figure 4** Visualization of the structures used to compute the MFE in the (a) general setting, (b) TR setting and (c) strand soup setting.

## 4     Interaction of Triplet Repeats

We now consider a set $R_0$ of triplet repeat strands. Our goal is to find the minimum free energy secondary structure for $R_0$. We defined the computational problem MFE STRAND INTERACTION in Section 2.2. In the base pair maximization model, this gives exactly the same definition as in [7], where the authors showed that the problem is **APX**-hard (and by that **NP**-hard) for the general (non-triplet) case. On the other hand, Dirks et al [9] gave a factorial-time algorithm for computing the partition function over multiple strands. In this section, we improve both results in the sense that on the one hand, we show that the problem is **NP**-hard in a reasonable energy model even if restricted to triplet repeats of one pattern, and on the other hand we give an exponential-time instead of factorial-time algorithm for the problem. However, notice that our exponential-time algorithm is designed for solving the MFE from an algorithmic perspective, as discussed in Section 2.2. If we want to give a penalty for rotational symmetries to the free energy of secondary structures, as described by Dirks et al. [9], the DP will not necessarily compute the MFE value. For the partition function, we can account for the algorithmic overcounting and additionally, if desired, for penalties for rotational symmetries.

### 4.1     General RNA-RNA interactions

The difficulty of the problem lies in the fact that we need to consider all possible circular permutations of strands. Instead of trying all of these circular permutations one by one and applying a classical single-stranded folding algorithm, we build up the values for all possible circular permutations while exploring all possible joint secondary structures. More specifically, we will consider structures consisting of a leftmost strand and its position, a rightmost strand and its position, as well as a set of strands which have to appear in between the leftmost and rightmost strand (without specifying the ordering of these strands).

We can formulate DP recurrences as follows: Let $E_{s_i, r_j}$ be the minimum free energy induced by the base pair between the $i$-th base of strand $s$ and the $j$-th base of strand $r$. In our DP equations, $R \subseteq R_0$ denotes the subset of still available strands, $s \in R$ the leftmost strand, $r \in R$ the rightmost strand, $1 \leq i \leq |s|$ the current position in $s$, $1 \leq j \leq |r|$ the current position in $r$, and $c \in \{0, 1, 2\}$ indicates whether $s$ and $r$ will be connected by a base pair (0: no base pair allowed, 1: at least one base pair required, 2: a base pair is not required; if the left and right strand are equal, then $c = 2$). The structures with which our algorithm works are visualized in Figure 4 (a). The main recurrences are as follows:

$$M_{R,s_i,r_j,c} = \min \begin{cases} \begin{cases} M_{R,s_{i+1},r_j,c} & \text{if } i+1 \leq |s| \\ \min_{t \in R, c' \in \{0,1\}} M_{R-\{s\},t_1,r_j,c'} - \mathbb{1}_{c'=0} K_{\text{assoc}} & \text{if } i+1 > |s| \text{ and } c \neq 1 \\ +\infty & \text{else} \end{cases} \\ \begin{cases} E_{s_i,r_j} + \bar{M}_{R,s_i,r_j,2} & \text{if } c \neq 0 \\ +\infty & \text{if } c = 0 \end{cases} \\ \min_{R',t,k} E_{s_i,t_k} + \bar{M}_{R',s_i,t_k,2} + \bar{M}_{(R-R')\cup\{s\},t_k,r_{j+1},c} \end{cases}$$

where

$$\bar{M}_{R,s_i,r_j,c} = \begin{cases} M_{R,s_{i+1},r_{j-1},c} & \text{if } i+1 \leq |s| \text{ and } j-1 \geq 1 \\ \min_{t \in R-\{s,r\},c' \in \{0,1\}} M_{R-\{s,r\},t_1,r_{j-1},c'} - \mathbb{1}_{c'=0} K_{\text{assoc}} & \text{if } i+1 > |s| \text{ and } j-1 \geq 1 \\ \min_{u \in R-\{s,r\},c' \in \{0,1\}} M_{R-\{s,r\},s_{i+1},u_{|u|},c'} - \mathbb{1}_{c'=0} K_{\text{assoc}} & \text{if } i+1 \leq |s| \text{ and } j-1 < 1 \\ \min_{t,u \in R-\{s,r\},c' \in \{0,1\}} M_{R-\{s,r\},t_1,u_{|u|},c'} - \mathbb{1}_{c'=0} K_{\text{assoc}} & \text{else} \end{cases}$$

and $-K_{\text{assoc}}$ is a reward for an additional complex. We give this reward each time we "choose" a new strand from $R$ and decide that it should not be connected to the other extremity of the interval ($c' = 0$). The $\bar{M}_{R,s_i,r_j,c}$ equation gives the MFE for the region $]s_i, r_j[$ (i.e. $[s_{i+1}, r_{j-1}]$ if $i+1 \leq |s|$ and $j-1 \geq 1$, and introducing new strands in the other cases). The minimization requires some more detailed conditions which can be found in the appendix.

Choosing an arbitrary strand $s$, the minimum free energy can be finally computed by

$$E^*(R) = (m-1) \cdot K_{\text{assoc}} + \min_{r \in R-\{s\}, c \in \{0,1\}} M_{R,s_1,r_{|r|},c}$$

and the optimal secondary structure can be obtained through backtracking.

For the initialization, we can set $M_{\{s\},s_i,s_j} = 0$ for valid indices $j - i \leq \theta$ for any $s \in R$. The correctness of the algorithm and its running time are proven in the appendix. With $n$ denoting the length of the concatenation of all strand sequences in $R$, we obtain:

▶ **Theorem 14.** *MFE STRAND INTERACTION can be solved in time $\mathcal{O}(3^m \cdot n^3)$.*

▶ Remark 15. The above DP scheme could be transformed to compute the partition function, using a change of algebra. The issues of overcounting and rotational symmetries raised by Dirks *et al* [9] may then be addressed by computing a separate PF $Q^{[r]} = \sum_{S \in \Omega: S \text{ is } r\text{-symm.}} e^{-E(S)/kT}$ for each possible symmetry value $r$ (while subtracting the values for higher symmetry structures), and obtain the PF by summing these values.

## 4.2 Strand interactions for triplet repeats

We now consider the special case where all strands in our pool are triplet repeats. We call this restricted problem MFE TRIPLET REPEAT STRAND INTERACTION. Assume first that all strands have the same pattern and that we have a bounded number of different strand-lengths $p := |\{i \mid \exists r \in R : |r| = i\}|$. Regardless of the ordering of the strands, the resulting sequence of the concatenated strands is identical. We can therefore focus on the length of the strands and disregard their actual sequence.

We do not iterate over all subsets of $R$, since we only need to distinguish the number of strands of a certain length in the subset, in a count-sort-like manner. Thus we can represent a subset $R' \subseteq R$ by $(a_1, ..., a_p)$ where $a_i := |\{r \in R' \mid |r| = n_i\}|$ is the number of strands of size $n_i$ in $R$. An example is given in Figure 4 (b). As argued in the appendix, the exponent only depends on $p$, and for $n := \max_{r \in R} |r|$ we get:

▶ **Theorem 16.** *There is an XP algorithm for MFE* TRIPLET REPEAT STRAND INTERACTION *parametrized by the number of different lengths $p$, running in $\mathcal{O}((\frac{m}{p})^{2p} \cdot n^3 \cdot p)$ time.*

Notice that this algorithm can be extended to the case where we have different triplet patterns; the parameter then becomes the number of non-identical strands.

## 4.3 Computational hardness

In this subsection, we show that the parametrized approach seen before is the best we can hope for, and that, even for triplet repeats, the problem of deciding whether there is a secondary structure for $R_0$ with a free energy below a certain threshold $t$ is **NP**-complete, for a reasonable energy model. Note that for the general (non-triplet) case, this has already been shown in [7]. Our result is surprising in the sense that the concatenation of TR strands always yields the same sequence, and the only additional difficulty compared to the single-stranded case arises from the fact that we do not know the indices of the strand borders.

Our reduction requires more than the naive base pair maximization model, but to keep the reduction simple, we will not use the full Turner energy model. Instead, each base pair gives a free energy reward of $E^{\text{bp}} = -\frac{m}{3}$, where $m > 0$ is the number of interacting strands, while subdividing an interval into two intervals that are not strand-disjoint gives a multiloop penalty of $K_{\text{multi}} = +1$. Furthermore, each connected component reduces the strand association penalty by $-K_{\text{assoc}} := -1$. Finally, every hairpin loop must enclose at least three unpaired bases ($\theta = 3$). This model is extendable to the Turner model by setting equal energy values for interior and hairpin loops and account for the multiloop penalty in the corresponding energy values.

Let us define the main decision problem:

---

TRIPLET REPEAT MULTI-STRAND MFE
**Input:** A set $R$ of explicitly encoded triplet repeat strands of the same pattern and a target free energy value $t$.
**Output:** Is there a secondary structure $S \in \Omega(R)$ with $E(R, S) \leq t$?

---

Even if the following reduction does not work in the base pair maximization model, a DP algorithm for base pair maximization in this setting seems unlikely, as, under the assumption $\mathbf{P} \neq \mathbf{NP}$, one would not be able to generalize the algorithm to more complex energy models.

We will show **NP**-hardness by reduction from the following problem:

---

SUMMING TRIPLES
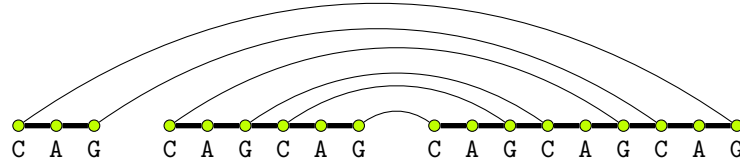**Input:** list of distinct positive integers $s_1, ..., s_{3n}$, encoded in unary
**Output:** Is there a partition of the input into triples $(a_i, b_i, c_i)$ such that $a_i + b_i = c_i$?

---

This has been shown to be strongly NP-hard [18]. We define $v := \sum_{i=1}^{3n} s_i$.

The reduction is as follows: We create a strand $r_i := (CAG)^{s_i}$ for each integer $s_i$. Hence, we have $n = \frac{m}{3} = -E^{\text{bp}}$. We denote by $R$ the set of strands. We set the target minimum free energy to $t := -(3v + 1)n$.

Assume that there is a partition into summing triples. Our secondary structure is built such that for each triple $a + b = c$, we add the base pairs

$$(a_1, c_{|c|}), (a_3, c_{|c|-2}), (a_4, c_{|c|-3}), (a_6, c_{|c|-5}), ..., (a_{|a|-2}, c_{|c|-|a|+3}), (a_{|a|}, c_{|c|-|a|+1}),$$
$$(b_1, c_{|c|-|a|}), (b_3, c_{|c|-|a|-2}), ..., (b_{|b|-2}, c_3), (b_{|b|}, c_1)$$

**Figure 5** Optimal secondary structure corresponding to a valid summing triple $(1, 2, 3)$. The $1 \cdot 3 + 2 \cdot 3$ bases on the left perfectly match to the $3 \cdot 3$ bases on the right.

Note that all base pairs are labeled with $C - G$ or $G - C$. Figure 5 visualizes the secondary structure for the exemplary triple $1 + 2 = 3$. We claim that $S$ is unpseudoknotted for the circular permutation $a_1 \cdot b_1 \cdot c_1 \ldots a_n \cdot b_n \cdot c_n$ and that $E(R, S) = t$.

Since any two triples of strands are not connected, we have exactly $n$ connected components. Each connected component consists of one large stacked loop with innermost base pair $(b_{|b|}, c_1)$ (i.e. we do not violate the constraint that every innermost base pair must include three unpaired bases, because the base pair is inter-strand). Since $a + b = c$, the outermost base pair is $(a_1, c_{|c|})$. There is no multiloop involved in $S$, so each triple $(a_i, b_i, c_i)$ contributes a free energy of $2|c| \cdot E^{\mathrm{bp}} - K_{\mathrm{assoc}} = -6n|c| - 1$. Since all triples are correctly summing, we have $\sum_{i=1}^{n} c_i = \frac{1}{2} v$. Thus indeed the minimum free energy is at most

$$\sum_{i=1}^{n} -6n|c_i| - 1 = -6n \sum_{i=1}^{n} |c_i| - n = -6n \cdot \frac{1}{2} v - n = -3nv - n = t$$

Before showing the opposite direction, we introduce the following simple lemmata:

▶ **Lemma 17.** *If some $C$ or $G$ base remains unpaired in a secondary structure $S$, $E(R, S) > t$.*

**Proof.** First notice that in every valid secondary structure, all $A$ bases remain unpaired (since there are no $U$ bases). There are $2v$ bases of $C/G$ in total. Since we assumed that one of them is unpaired, there can be at most $v - 1$ base pairs. We can have at most $3n$ complexes, so the strand association penalty is reduced by at most $3n$. Thus we have $E(R, S) \geq -3n(v - 1) - 3n = -3vn > -(3v + 1)n = t$.     ◀

▶ **Lemma 18.** *If $S$ contains a hairpin loop, $E(R, S) > t$.*

**Proof.** A hairpin loop must enclose at least three bases. Since in the `CAG` triplet pattern any two consecutive bases involve at least one $C$ or $G$, we can apply Lemma 17 and conclude.     ◀

Now assume for an arbitrary $S \in \Omega$ that $E(R, S) \leq t$. We first show that there must be exactly $n$ connected components, each with three strands. Assume that there is a connected component with less than three strands. If it has only one strand, it must contain a hairpin loop, and by Lemma 18, $E(R, S) > t$. If the complex contains two strands, first of all the two strands have a different number of triplet repeats, since all $s_i$ are distinct. This implies that if the innermost loop is inter-strand (if it is intra-strand we again apply Lemma 18) and has no multiloop, some $G$ or $C$ base must be unpaired (since base pairs can then only be between the two strands, but one of the strands contains at least one $G$ and one $C$ base more than the other). Then, by Lemma 17, $E(R, S) > t$. If it has a multiloop, there have to be two innermost base pairs, one of which must be intra-strand, and we can apply Lemma 18.

Since we ruled out complexes of one or two strands and the total number of strand is divisible by 3, we know that if there is a complex with four strands, our secondary structure will have $< n$ connected components. Thus the best achievable score will be $-n+1-3nv > t$. Hence, any $S \in \Omega$ with $E(R, S) \leq t$ consists of $n$ complexes, each consisting of three strands $a_i, b_i, c_i$ with $|a_i| < |b_i| < |c_i|$. We claim that for all $i \in [n]$, $|a_i| + |b_i| = |c_i|$.

By contradiction, assume $|a_i| + |b_i| \neq |c_i|$ and first consider the case that there are no multiloops. This implies that there is only one innermost base pair. If it is intra-strand, we obtain a contradiction to $E(R, S) \leq t$ by Lemma 18. If it is inter-strand, all remaining base pairs must be between one of two strands $d, e$ on the one side and the third strand $f$ on the other side. Since $|d| + |e| \neq |f|$ for any such partition, one of the two sides will be left with at least one unpaired $G$ and one unpaired $C$, and we apply Lemma 17.

Now we consider the case of multiloops. Any multiloop where the cutpoint between the two recursive structures is on a strand border (and thus is not penalized) implies an innermost base pair in both recursive structures, and since by pigeonhole principle one of the two recursive structures is single-stranded, we have a hairpin loop and $E(R, S) > t$ by Lemma 18. In the other case, we have a multiloop penalty of $+1$. Thus we can lower bound $E(R, S) \geq -n - 3nv + 1 > t$.

This finishes the proof that $|a_i| + |b_i| = |c_i|$, and we get $\frac{|a_i|}{3} + \frac{|b_i|}{3} = \frac{|c_i|}{3}$. By the construction, each strand $r$ corresponds to one integer $\frac{|r|}{3}$ in the set of integers of our original instance. Thus, $(\frac{|a_i|}{3}, \frac{|b_i|}{3}, \frac{|c_i|}{3})$ for all complexes $\{a_i, b_i, c_i\}$ for $1 \leq i \leq n$ is a valid set of summing triples.

The reduction is polynomial-time, since in the SUMMING TRIPLES problem, the integers are encoded in unary. Membership in **NP** follows by the fact that we can evaluate the energy given a secondary structure and its unpseudoknotted circular permutation.

▶ **Theorem 19.** *UNARY TRIPLET REPEAT MULTI-STRAND MFE is **NP**-complete.*

## 4.4 Strand soup interaction

We now consider the computational problem MFE STRAND SOUP INTERACTION as defined in Section 2.2. We can adapt the algorithm from above and we do not need to keep track of the (exponentially many) subsets anymore, yielding a polynomial-time algorithm. We do not charge any strand association penalty, since we require one single complex anyways. However, we still must enforce connectivity. To this end, we encode by $c = 1$ that $s$ and $r$ still need to be connected, and by $c = 2$ that they already are connected. Furthermore, instead of keeping track of a subset of remaining strands, we just need the number of remaining strands $m$, as seen in Figure 4 (c). We obtain the following DP equations:

$$
M_{m,s_i,r_j,c} = \min \begin{cases} \begin{cases} M_{m,s_{i+1},r_j,c} & \text{if } i+1 \leq |s| \\ \min_{t \in R} M_{m-1,t_1,r_j,1} & \text{if } i+1 > |s| \text{ and } c \neq 1 \\ +\infty & \text{else} \end{cases} \\ E_{s_i,r_j} + \bar{M}_{m,s_i,r_j,2} \\ \min_{m',t,k \text{ s.t. } (\ast)} E_{s_i,t_k} + \bar{M}_{m',s_i,t_k,2} + \bar{M}_{m-m'+1,t_k,r_{j+1},c} \end{cases}
$$

where

$$
\bar{M}_{m,s_i,r_j,c} = \begin{cases} M_{m,s_{i+1},r_{j-1},c} & \text{if } i+1 \leq |s| \text{ and } j-1 \geq 1 \\ \min_{t \in R} M_{m-1,t_1,r_{j-1},1} & \text{if } i+1 > |s| \text{ and } j-1 \geq 1 \\ \min_{u \in R} M_{m-1,s_{i+1},u_{|u|},1} & \text{if } i+1 \leq |r| \text{ and } j-1 < 1 \\ +\infty & \text{else} \end{cases}
$$

The minimum free energy can be finally computed by

$$E^*(R, m) = \min_{s,r \in R} M_{m,s_1,r_{|r|},1}$$

and the optimal secondary structure can be obtained through backtracking. We initialize $M_{1,s_i,s_j,2} = 0$ for all $j - i \leq \theta$.

The correctness mostly follows from Section 4.1, but we still have to argue that we correctly minimize over *connected* secondary structures only, which is done in the appendix.

Regarding the running time, the table size is bounded by $m \cdot p^2 \cdot n^2 \cdot 3$, where $n := \max_{s \in R} |s|$. The running time to compute one table entry is dominated by the last case, where we minimize over $\mathcal{O}(m \cdot p \cdot n)$ cutpoints and need $\mathcal{O}(p)$ time for each new strand. In total, we obtain an algorithm with running time $\mathcal{O}(n^3 \cdot m^2 \cdot p^4)$. We can then conclude:

▶ **Theorem 20.** *MFE* Unlimited Strand Interaction *can be solved in time $\mathcal{O}(n^3 \cdot m^2 \cdot p^4)$.*

▶ Remark 21. Additionally to restricting the number of interacting strands, one can extend the above algorithm to restrict the size of the concatenated sequence. This is possible by keeping track of the current size of the sub-interval in the DP tables, and updating these values whenever a new strand is introduced.

This might be useful if the sequences in the base set have different length, as the basic algorithm would favor larger sequences because they usually allow for more base pairs.
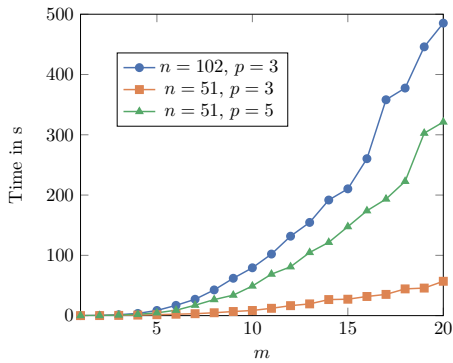
▶ Remark 22. The case of triplet repeats gives a slight improvement to the running time. Since all strands look the same except for their length, we can use a table with entries of the form $M_{m,i,j,c}$, where $i$ and $j$ denote the remaining number of bases in the leftmost and rightmost strand. This reduces the space complexity to $\mathcal{O}(m \cdot n^2)$, but the computation of one table entry still takes the same time, giving an overall time complexity of $\mathcal{O}(n^3 \cdot m^2 \cdot p^2)$.
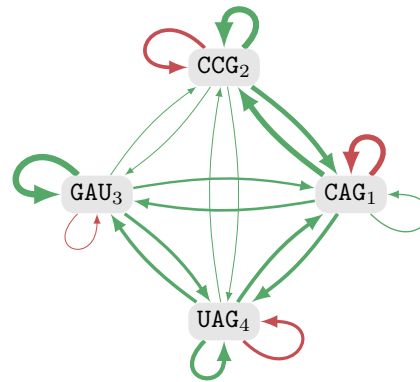
## 5   Empirical proof of concept

The goal of this section is to show how the algorithms described in the previous section can be used to answer biologically relevant questions regarding triplet repeats. We implemented the algorithm described in Section 4.4, which hereunder we call SoupFold, as well as its partition function equivalent, together with a (stochastic) backtracking procedure. Since we only limit the number of interacting strands but not their size, without further restrictions, the program would prefer large strands since they usually give more base pairs. To counteract this effect, we introduce a penalty on the length of a strand. Note that one could also set a maximum length of the concatenated sequence, as described in Remark 21. The empirically observed running time matches the theoretical running time well, as can be seen in Figure 6. The source code to reproduce analyses is available at:

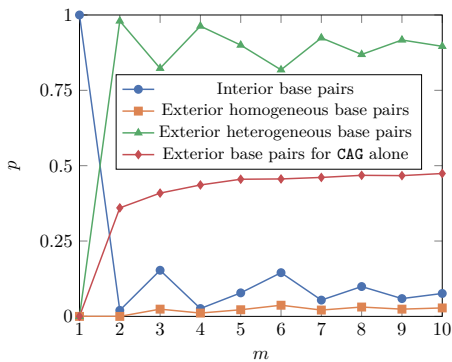https://github.com/kimonboehmer/soupfold/

Regarding the stochastic backtracking, we must account for the overcounting of rotationally asymmetric secondary structures as well as for the overcounting because of the positioning of different connected components. We address these two issues by rejection sampling. In theory, it is also necessary to adjust the overcounting correction for rotationally symmetric structures (because they are overcounted less often) but our experiments showed that the observed probability of encountering such rotational symmetries is 0 for triplets with 15 repeats or more. Thus, for efficiency reasons, we do not include this case in our rejection sampling, arguing that the changes to the probability would be too small to observe.
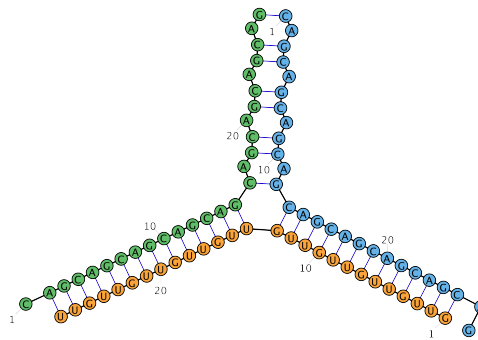
**Figure 6** Empirical running time for increasing $m$ and various values for $n$ and $p$.



**Figure 7** Affinity of triplet repeats (%BPs in soup model, coded by line thickness) for external (green) and internal (red) interactions.



**Figure 8** Probability $p$ that a certain type of base pair is observed for increasing #strands $m$, either in a soup $\{\texttt{CAU}_{20}, \texttt{GGG}_{20}\}$, or for $\texttt{CAG}_{20}$.



**Figure 9** Exemplary MFE structure for strand pool $\{(\texttt{GUU})^9, (\texttt{CAG})^9, (\texttt{ACG})^9\}$ computed by Soup-Fold with $m = 3$ (RiboSketch [16]).

## 5.1 Homogeneous triplet soup

We first consider the case where all strands are of the same pattern. The MFE of a soup of homogenous triplets behaves canonically, in the sense that all folding patterns have almost identical MFE structures (as can be expected, considering our results on single-strand TR in Section 3). Furthermore, we observed that the number of base pairs increases canonically with the sequence length and with the number of interacting strands (except for the case of only one strand, where we loose one base pair due to a hairpin loop).

## 5.2 Heterogeneous triplet soup

More interesting observations can be made in a heterogeneous pool. We can observe that different TR pattern strands can achieve more base pairs than the theoretical upper bound for a homogeneous strand pool (see Figure 9).

In order to assess the capability of different strand soups to form droplets, we want to determine the probability of a base pair in the Boltzmann ensemble being between two strands (exterior) as opposed to folding (interior). If the strand soup consists only of triplets of one pattern, all exterior base pairs will be homogeneous, as opposed to heterogeneous for

an interaction of two strands of different patterns. In the homogeneous case, we can observe an increase of exterior base pairs for increasing number of interacting strands $m$, as presented by the red line in Figure 8. The probabilities in a setting with strands of different patterns are much richer and less canonical, as can be seen at the example of the interaction of `CAU` and `GGG`, presented by the other lines in Figure 8. These probabilities highly depend on the number of strands, and only start to "converge" with quite high values of $m$.

To obtain a broader picture, we performed stochastic backtracking on all possible $4^6$ pairs of triplet repeat patterns $\{TVW, XYZ\}$ as strand sets, with $m$ between 2 and 5, and computed the probability of a base pair being interior, exterior-homogeneous or exterior-heterogeneous. A visualization and a small discussion can be found in the appendix. From a synthetic biology perspective, some triplet repeats aggregate and form a Liquid-Liquid Phase Separation, which can be used to isolate subprocesses, thereby implementing a notion of orthogonality. In order to maximize the number of independent tasks being performed by a modified bacteria, it would then be desirable to find a large number of triplet repeat patterns such that the probability of heterogeneous base pairs is low.

For that, we can model the patterns as vertices of a graph and draw an edge if the heterogeneous base pair probability between two patterns for $m = 5$ is high (we set the threshold to 0.175). We then want to determine a maximum independent set (MIS), i.e. the largest set of triplets that do not have a high probability of interacting pairwise with each other. We used an exact solver [11] to obtain a MIS of size 4, namely `CAG, CCG, GAU, UAG`.

We then executed our algorithm on these triplet patterns as strand soup, and could indeed observe that the probability of exterior heterogeneous base pairs is clearly below 0.2 for values of $m$ between 1 and 10. In Figure 7, we depict the number of base pairs that are between two types of strands, for $m = 5$ and our four independent TR patterns as strand soup. We added a bonus to the appearance of strands, to ensure that all strands of the soup appear equally often in the constructed structures. We observe that for three of the four triplets, for exterior base pairs, the most likely interacting strand is of the same type.

## 6    Conclusion and Discussion

In this work, we investigated the algorithmic aspects of folding and interactions of triplet repeat RNA sequences, while also revisiting the general (non-triplet) setting in the interaction setting. For the folding of individual triplets, we found that their repetitive structure allows us to immediately characterize the MFE and partition function value, without the need of a more time-consuming DP approach. For interactions of RNA sequences, we exhibited a new algorithm with improved running time that avoids the factorial-time iteration over all permutations and acts as a foundation for the design of specialized algorithms, as the XP algorithm for triplet repeats. For the "strand soup" setting, we derived a polynomial-time algorithm and demonstrated possible uses for experiments regarding triplet repeats.

For future work, it is desirable to describe in detail how to extend the MFE STRAND INTERACTION algorithm to the full thermodynamic setting considered in [9]. While the extension to the Turner model does not pose any algorithmic challenges, it would be interesting to implement a variant of the inside/outside algorithm to compute exactly base-pairing probabilities and other expected values of additive properties. Finally, the joint conformation space explored in this work is heavily restricted by the existence of a non-crossing strand ordering. More complex conformational spaces could be captured by using DP approaches akin to those used to include pseudoknots in RNA structure prediction.

────────── **References** ──────────

**1**  Dilimulati Aierken and Jerelle A Joseph. Accelerated simulations of rna phase separation: a systematic study of non-redundant tandem repeats. *bioRxiv*, pages 2023–12, 2023.

**2**  Can Alkan, Emre Karakoc, Joseph H Nadeau, S Cenk Sahinalp, and Kaizhong Zhang. Rna–rna interaction prediction and antisense rna target search. *Journal of Computational Biology*, 13(2):267–282, 2006.

**3**  Kimon Boehmer, Sarah J. Berkemer, Sebastian Will, and Yann Ponty. Soupfold. Software, version 1.0., ANR-funded SYNORG project (ANR-23-CE44-0027), swhId: `swh:1:dir:3cb0c3e63356a53b4fa150d14e4c273678ef638d` (visited on 2024-08-16). URL: `https://github.com/kimonboehmer/soupfold/`.

**4**  Alin Bostan, Frédéric Chyzak, Gr égoire Lecerf, Bruno Salvy, and Éric Schost. Differential equations for algebraic functions. In C. W. Brown, editor, *ISSAC'07: Proceedings of the 2007 international symposium on Symbolic and algebraic computation*, pages 25–32. ACM Press, 2007. `doi:10.1145/1277548.1277553`.

**5**  Karl Bringmann, Fabrizio Grandoni, Barna Saha, and Virginia Vassilevska Williams. Truly subcubic algorithms for language edit distance and rna folding via fast bounded-difference min-plus product. *SIAM Journal on Computing*, 48(2):481–512, 2019. `doi:10.1137/17M112720X`.

**6**  Yi-Jun Chang. Hardness of rna folding problem with four symbols. *Theoretical Computer Science*, 757:11–26, 2019. `doi:10.1016/j.tcs.2018.07.010`.

**7**  Anne Condon, Monir Hajiaghayi, and Chris Thachuk. Predicting minimum free energy structures of multi-stranded nucleic acid complexes is apx-hard. In *27th International Conference on DNA Computing and Molecular Programming (DNA 27)(2021)*. Schloss-Dagstuhl-Leibniz Zentrum für Informatik, 2021.

**8**  A. Denise, Y. Ponty, and M. Termier. Controlled non-uniform random generation of decomposable structures. *Theoretical Computer Science*, 411(40):3527–3552, 2010. `doi:10.1016/j.tcs.2010.05.010`.

**9**  Robert M Dirks, Justin S Bois, Joseph M Schaeffer, Erik Winfree, and Niles A Pierce. Thermodynamic analysis of interacting nucleic acid strands. *SIAM review*, 49(1):65–88, 2007.

**10**  Haotian Guo, Joseph C Ryan, Xiaohu Song, Adeline Mallet, Mengmeng Zhang, Victor Pabst, Antoine L Decrulle, Paulina Ejsmont, Edwin H Wintermute, and Ariel B Lindner. Spatial engineering of E. coli with addressable phase-separated RNAs. *Cell*, 185(20):3823–3837, 2022.

**11**  Fanny Hauser, Ferdinand Ermel, and Kimon Boehmer. Clique cover based vertex cover solver. `https://github.com/f-erm/CliqueCoverBasedVertexCoverSolver`, 2024.

**12**  Liang Huang, He Zhang, Dezhong Deng, Kai Zhao, Kaibo Liu, David A Hendrix, and David H Mathews. LinearFold: linear-time approximate RNA folding by 5'-to-3' dynamic programming and beam search. *Bioinformatics*, 35(14):i295–i304, July 2019. `doi:10.1093/bioinformatics/btz375`.

**13**  Atagun U Isiktas, Aziz Eshov, Suzhou Yang, and Junjie U Guo. Systematic generation and imaging of tandem repeats reveal base-pairing properties that promote RNA aggregation. *Cell Reports Methods*, 2(11), 2022.

**14**  Ryo Kurokawa, Mariko Kurokawa, Akihiko Mitsutake, Moto Nakaya, Akira Baba, Yasuhiro Nakata, Toshio Moritani, and Osamu Abe. Clinical and neuroimaging review of triplet repeat diseases. *Japanese Journal of Radiology*, 41(2):115–130, 2023.

**15**  L. Lipshitz. *D*-finite power series. *Journal of Algebra*, 122(2):353–373, 1989.

**16**  Jacob S Lu, Eckart Bindewald, Wojciech K Kasprzak, and Bruce A Shapiro. RiboSketch: versatile visualization of multi-stranded RNA and DNA secondary structure. *Bioinformatics*, 34(24):4297–4299, June 2018. `doi:10.1093/bioinformatics/bty468`.

**17**  Hiranmay Maity, Hung T Nguyen, Naoto Hori, and D Thirumalai. Odd–even disparity in the population of slipped hairpins in rna repeat sequences with implications for phase separation. *Proceedings of the National Academy of Sciences*, 120(24):e2301409120, 2023.

**18**  Colin McDiarmid. Pattern minimisation in cutting stock problems. *Discrete applied mathematics*, 98(1-2):121–130, 1999.

**19**   R Nussinov and A B Jacobson. Fast algorithm for predicting the secondary structure of single-stranded rna. *Proceedings of the National Academy of Sciences*, 77(11):6309–6313, 1980. `doi:10.1073/pnas.77.11.6309`.

**20**   B. Salvy and P. Zimmerman. GFUN: a Maple package for the manipulation of generating and holonomic functions in one variable. *ACM Transactions on Mathematical Software*, 20(2):163–177, 1994.

**21**   Sharan R. Srinivasan, Claudio Melo de Gusmao, Joanna A. Korecka, and Vikram Khurana. Chapter 18 - repeat expansion disorders. In Michael J. Zigmond, Clayton A. Wiley, and Marie-Francoise Chesselet, editors, *Neurobiology of Brain Disorders (Second Edition)*, pages 293–312. Academic Press, second edition edition, 2023. `doi:10.1016/B978-0-323-85654-6.00048-4`.

**22**   Douglas H. Turner and David H. Mathews. NNDB: the nearest neighbor parameter database for predicting stability of nucleic acid secondary structure. *Nucleic Acids Research*, 38(suppl_1):D280–D282, October 2009. `doi:10.1093/nar/gkp892`.

**23**   Michael Zuker and Patrick Stiegler. Optimal computer folding of large RNA sequences using thermodynamics and auxiliary information. *Nucleic Acids Research*, 9(1):133–148, January 1981. `doi:10.1093/nar/9.1.133`.

## A   Appendix for Section 3

## A.1   Proof for Lemma 6

▶ **Lemma 6.** *Let $\theta > 1$. The minimum MFE structure of a folding sequence $(XYZ)^k$ has value*

- $k - 1 - \frac{\theta-1}{3}$, *if* $(\{X,Z\} \notin P \wedge (\theta + 3k) \equiv_6 4) \vee (\{X,Y\}, \{Y,Z\} \notin P \wedge (\theta + 3k) \equiv_6 1)$
- $k - \lfloor \frac{\theta+1}{3} \rfloor$, *otherwise*

*Furthermore, a minimum MFE structure is obtained by a single helix of base pairs of one letter pair $p$. If both $\{X,Z\} \in P$ and one of $\{X,Y\}$ and $\{Y,Z\} \in P$, we set $p := \{X,Z\}$ if $(\theta + 3k) \equiv_6 4$ and $p := \{X,Y\}$ (or $p := \{Y,Z\}$) if $(\theta + 3k) \equiv_6 1$; otherwise, we set $p$ to the letters of an arbitrary pairable base pair.*

**Proof.** We start by showing that the corresponding secondary structures achieve the claimed score. By Observation 2, we only need to consider $\theta \equiv_3 0$ and $\theta \equiv_3 1$.

First assume $\{X,Z\} \in P$ and $\{X,Y\}, \{Y,Z\} \notin P$. We will derive the other cases from this one. Consider a large stacking of $X - Z$ bases. If $\theta = 3$, we only cannot match the $X - Z$ pair of the innermost repeat in the case $k \equiv_2 1$ and we only cannot match the $Z - X$ pair between the two innermost repeats in the case $k \equiv_2 0$. For all other pairs of repeats we obtain exactly two base pairs and hence we get $k - 1 = k - \lfloor \frac{\theta+1}{3} \rfloor$ base pairs. Inductively, let us show that we can obtain $k - \lfloor \frac{\theta'+1}{3} \rfloor$ base pairs for $\theta' := \theta + 3$. In other words, we only need to show that by increasing $\theta$ by 3, we get one base pair less. If the innermost base pair is $X - Z$, its enclosed region starts and ends with a $Y$ and there are currently at least $\theta + 1$ free enclosed bases (because the region is of the form $Y(ZXY)^{\theta/3}$), and by deleting the $X - Z$ base pair, we obtain $XY(ZXY)^{\theta/3}Z$, that is $\theta + 3$ enclosed bases. Else, for a $Z - X$ base pair, the region has the form $(XYZ)^{\theta/3}$. After deleting the innermost base pair $Z - X$, the new enclosed region starts and ends with a $Y$ (the region is of the form $YZ(XYZ)^{\theta/3}XY$), so there are at least $\theta + 4$ enclosed bases. Thus we can achieve $k - \lfloor \frac{\theta+1}{3} \rfloor$ base pairs.

If $\theta \equiv_3 1$, we distinguish two equivalence classes: In the first, $k$ is even and $\theta \equiv_6 1$ *or* $k$ is uneven and $\theta \equiv_6 4$, and in the second equivalence class, we have the other two cases.

For $\theta = 4$, for $k \equiv_2 1$, our lemma only claims $k - 2$ base pairs. We can indeed leave the innermost repeat as well as the next $Z - X$ pair unpaired, and greedily create stackings outside of this region, obtaining $k - 2$ base pairs. For $k \equiv_2 0$, We can proceed as for the even case in $\theta = 3$.

Consider $\theta + 3$ now. We add an unpaired triplet in the middle of the sequence. Now, the number of base pairs is equal to the case $k - 1$ (of opposite parity) with $\theta$ enclosed bases.

We thus established the lower bound for the $\{X, Z\} \in P$ case. For the "otherwise"-case, Lemma 4 already gives us the required upper bound. Therefore, we only need to argue about the upper bound $k - 1 - \frac{\theta - 1}{3}$ in the case that $\{X, Y\}, \{Y, Z\} \notin P$ and $(\theta + 3k) \equiv_6 1$. Assume a secondary structure that achieves more base pairs. Firstly, we cannot have any multiloops or exterior loops since that would imply two regions of unpaired enclosed bases, which then only allows $k - 2\lfloor\frac{\theta+1}{3}\rfloor \leq k - 1 - \frac{\theta-1}{3}$ base pairs. Additionally, for each secondary structure $S$ with $i < j'$ and $k > 0$ such that $\{i, j'\} \in S$ and the interval $[j' + 1, j' + 3k]$ only consists of unpaired bases, we can delete the base pair $\{i, j'\}$ and instead add base pair $\{i, j' + 3k\}$ without reducing the number of base pairs. In other words, for any interval, it is always better to pair the leftmost base to the rightmost possible base than to any other interior base. We thus only need to consider the canonical structures of $X - Z/Z - X$-stackings.

Consider an odd $k$ with all base pairs in the canonical way (for $\theta = 4$). The innermost triplet repeat bases $X$ and $Z$ have to stay unpaired, as well as the $Z$ and $X$ which are adjacent to that repeat. The innermost base pair $X - Z$ now has $7 = \theta + 3$ enclosed bases. We thus have $k - 2$ base pairs. Inductively, for $\theta' := \theta + 6$, the next two innermost base pairs will have $\theta + 3 < \theta'$ and $\theta + 3 + 2 < \theta'$ enclosed bases, thus are both not available.

Consider an even $k$ with all base pairs in the canonical way (for $\theta = 7$). The two innermost triplet repeats have to stay unpaired, as well as the $Z$ and $X$ which are adjacent to that repeat. The innermost base pair $X - Z$ now has $10 = \theta + 3$ enclosed bases. The rest of the argument is exactly as above.

If $\{X, Z\} \notin P$, we can assume without loss of generality that $\{X, Y\} \in P$ (the arguments are symmetrical for $\{Y, Z\} \in P$, and we assumed to have a folding strand). We can reduce any such instance $(XYZ)^k$ to $(YZX)^{k-1}$ (by letting out the leftmost $X$ and the rightmost $Y$ and $Z$, and implicitly pairing these outermost $X$ and $Y$, which is always optimal). Thus, all results can be directly obtained from the case $\{X, Z\} \in P$, by changing odd and even. The upper bound can also be derived by that.                                                                                                          ◀

## B    Appendix for Section 4

### B.1    Proof of correctness for the exponential-time algorithm

We now prove that $M_{R,s_i,r_j}$ is computed correctly. By slight abuse of notation, we write $s_i \in S$ for $s_i \in \bigcup_{P \in S} P$.

▶ **Definition 23.** *An interval for this DP is denoted by $[R, s_i, r_j, c]$ where $s, r \in R$, $1 \leq i \leq |s|$, $1 \leq j \leq |r|$ and $c \in \{0, 1, 2\}$. An interval $[R', t_k, u_\ell, c']$ is **included** in interval $[R, s_i, r_j, c]$, written $[R', t_k, u_\ell, c'] \preccurlyeq [R, s_i, r_j, c]$, if one of the following holds:*

- $R' \subset R$ *and* $|R'| < |R| - 1$
- $R' \subset R$, $|R'| = |R| - 1$ *and* $s = t \vee r = u$
- $R' = R$, $s = t$, $r = u$, $i \leq k$ *and* $\ell \leq j$.

*If we replace* both *inequalities by strict inequalities in the last point, the interval is **strictly included** and we write $[R', t_k, u_\ell, c] \prec [R, s_i, r_j, c]$.*

Each such interval is associated to a minimum free energy as follows:

▶ **Definition 24.** *Let $I := [R, s_i, r_j, c]$. $\Omega(I)$ is the set of all secondary structures that are valid for this interval, or more formally, a secondary structure $S$ must fulfill:*

- $S \in \Omega(R)$
- $s_k, r_\ell \notin S$ for any $k < i$ and $\ell > j$
- $c = 1$ implies the existance of a base pair between $s$ and $r$ (that is, $\{s_k, r_\ell\} \in S$ for some $i \leq k \leq |s|, 1 \leq \ell \leq j$) and $c = 0$ implies that there is no such base pair.

The minimum free energy of $I$ is defined as $MFE(I) := \min_{S \in \Omega(I)} E(R, S)$.

The minimum free energy of an open interval $MFE(]R, s_i, r_j, c[)$ is the minimum free energy over all secondary structures and all intervals $I' \prec I$ where $c$ specifies the connectedness of $s$ and $r$.

We also observe that an optimal structure is optimal for any substructure that includes all its base pairs:

▶ **Observation 25.** *If $E(R, S) = MFE([R, s_i, r_j, c])$ and $S$ only contains base pairs in some interval $[R', t_k, u_\ell, c] \preccurlyeq [R, s_i, r_j, c]$, then $S = MFE([R', t_k, u_\ell, c])$.*

We first show that our helper equation $\bar{M}$ is computed correctly:

▶ **Lemma 26.** *Assuming that $M_{R', t_k, u_\ell, c'} = MFE(I' := [R', t_k, u\ell, c'])$ for all $I' \preccurlyeq I := [R, s_i, r_j, c]$, we have $\bar{M}_{R, s_i, r_j, c} = MFE(]R, s_i, r_j, c[)$.*

**Proof.** We distinguish four cases:
- **Case 1:** $i + 1 \leq |s|$ and $j - 1 \geq 1$. In that case, for any $I' \prec I$, we have $I' \preccurlyeq [R, s_{i+1}, r_{j-1}, c]$ and thus $MFE(I') \geq MFE([R, s_{i+1}, r_{j-1}, c]) = \bar{M}_{R, s_i, r_j, c}$ by assumption. Thus $MFE(]R, s_i, r_j, c[) = \bar{M}_{R, s_i, r_j, c}$.
- **Case 2:** $i + 1 > |s|$ and $j - 1 \geq 1$. For any $I' \prec I$, there is a $t \in R - \{s\}$ and a $c' \in \{0, 1\}$ with $I' \preccurlyeq [R - \{s\}, t_1, r_{j-1}, c']$. It thus suffices to minimize over the strands $R - \{s, r\}$ while taking into account a possible strand disconnection reward. We have $\min_{t \in R - \{s, r\}, c' \in \{0, 1\}} M_{R - \{s\}, t_1, r_{j-1}, c'} - \mathbb{1}_{c'=0} K_{\text{assoc}} = MFE(]R, s_i, r_j, c[)$.
- **Case 3:** $i + 1 \leq |s|$ and $j - 1 < 1$. This case is completely symmetrical to Case 2.
- **Case 4:** $i + 1 > |s|$ and $j - 1 < 1$. For any $I' \prec I$, there are $t, u \in R - \{s, r\}$ with $I' \preccurlyeq [R - \{s, r\}, t_1, u_{|u|}, 2]$. It thus suffices to minimize twice over the strands $R - \{s, r\}$ while taking into account a possible strand disconnection reward. We have $\min_{t, u \in R - \{s, r\}, c' \in \{0, 1\}} M_{R - \{s, r\}, t_1, u_{|u|}, c'} - \mathbb{1}_{c'=0} K_{\text{assoc}} = MFE(]R, s_i, r_j, c[)$.   ◄

▶ **Lemma 27.** *The algorithm computes the table entries correctly, i.e. $M_{R, s_i, r_j, c} = MFE([R, s_i, r_j, c])$ for all $R \subseteq R_0$, $s_i, r_j \in R$ and $c \in \{0, 1, 2\}$.*

**Proof.** We proceed by induction over the well-founded relation $\preccurlyeq$. Regarding the initialization, clearly no base pair can exist over an empty strand set, as well as over one strand where the number of enclosed base pairs between $i$ and $j$ is less than $\theta$. Therefore, these table entries are correctly initialized by 0.

Let us assume that all $M_{R', t_k, u_\ell, c}$ with $[R', t_k, u_\ell, c] \preccurlyeq [R, s_i, r_j, c]$ except $M_{R, s_i, r_j, c}$ itself have been computed correctly.
- **Case 1:** $s_i \notin S$. If $i + 1 \leq |s|$, we have $E(R, S) = MFE([R, s_{i+1}, r_j, c]) = M_{R, s_{i+1}, r_j, c}$ by Observation 25 and our induction hypothesis.
  Else, we first assume $c \neq 1$. Consider the strand $t$ that follows $s$ in the polymer graph of $S$ and consider the value $c'$ that specifies connectivity between $t$ and $r$ in $S$. Since $i$ is unpaired, we again have $E(R, S) = MFE([R - \{s\}, t_1, r_j, c']) - \mathbb{1}_{c'=0} K_{\text{assoc}} = M_{R - \{s\}, t_1, r_j, c} - \mathbb{1}_{c'=0} K_{\text{assoc}}$ as above.
  Finally, if $c = 1$, we look for the MFE of a structure in $[R, s_i, r_j, c]$ where $s$ and $r$ are connected by a base pair. Since there is only one base in $s$ remaining and we leave it unpaired, there is no such structure and thus $MFE([R, s_i, r_j, 1]) = +\infty$.

- **Case 2:** $S = \{\{s_i, r_j\}\} \cup S'$, where $S'$ is the best structure for any $I' \prec [R, s_i, r_j, c]$ with $s$ and $r$ arbitrarily connected (that is, $]R, s_i, r_j, 2[$). First assume $c \neq 0$. In this case, we have $E(R, S) = E_{s_i, r_j} + \text{MFE}(]R, s_i, r_j, 2[) = E_{s_i, r_j} + \bar{M}_{R, s_i, r_j, 2}$, where we could apply Lemma 26 because of the induction hypothesis.

  Now assume $c = 0$. We minimize over all structures such that $s$ and $r$ are not connected, but require $\{s_i, r_j\} \in S$. Thus $\text{MFE}([R, s_i, r_j, 0]) = +\infty$.

- **Case 3:** $S = \{s_i, t_k\} \cup S' \cup S''$ for some $t_k \neq r_j$, where $S'$ (resp. $S''$) is the best structure for any $I' \prec [R', s_i, t_k, c]$ (resp. $I' \prec [R'', t_k, r_j, c]$), with $R'$ being all strands between $s$ and $t$ in the polymer graph of $S$, and $R''$ being all strands between $t$ and $r$.

  Note that $s$ and $t$ are connected, thus in $S'$ the connectivity bit will be set to 2. On the other hand, the connectedness of $t$ and $r$ (for structure $S''$) is by transitivity of connectivity determined by the connectedness between $s$ and $r$, that is, $c$. We then have $\text{MFE}([R, s_i, r_j, c]) = E_{s_i, t_k} + \text{MFE}(]R', s_i, t_k, 2[) + \text{MFE}(]R'', t_k, r_j, c[)$.  ◄

We now briefly discuss the running time. The number of table entries is bounded by $2^m \cdot n^2$, where $n := \sum_{r \in R} |r|$ is the size of the concatenated sequence. The last case of the DP equation dominates the running time for computing one entry. In the worst case, we iterate over $2^{|R|}$ subsets and $n$ entries, which gives $\mathcal{O}(2^{|R|} \cdot n)$. Partitioned by subset size, we get

$$\sum_{t=0}^{m} \binom{m}{t} n^2 \cdot 2^t n = n^3 \cdot \sum_{t=0}^{m} \binom{m}{t} 2^t = n^3 \cdot \sum_{t=0}^{m} \binom{m}{t} 1^{m-t} 2^t = n^3 \cdot (1+2)^m = 3^m \cdot n^3$$

which bounds the total running time. Together with Lemma 27, we conclude.

**Detailled conditions and edge cases.**   When we minimize over all subsets, the following conditions must be respected:

$$\{s, t\} \subseteq R' \subseteq R \wedge 1 \leq k \leq |t| \wedge (k = |t| \rightarrow c \neq 1)$$
$$\wedge (s = t \rightarrow (k > i + \theta \wedge R' = \{s\}))$$
$$\wedge (r \in R' \rightarrow (t = r \wedge k < j \wedge R' = R \wedge c \neq 0))$$

We minimize over all possible triples $(R', t, k)$. A set $R'$ must clearly include $s$ and $t$ to form a valid interval and $k$ must be a valid position of $t$. If $s_i$ is paired to $t_{|t|}$, $s$ and $j$ are disconnected ($c \neq 1$). If $s = t$, we must respect $\theta$ and there is only one strand in $R'$. Finally, $r \in R'$ implies that $s_i$ forms a base pair with some base of $r$ (thus $t = r$ and $R' = R$), connectivity has to be allowed ($c \neq 0$) and $t_k$ must be in the interval ($k < j$). These conditions are sufficient and match our algorithm.

When we minimize over two new inner strands (in the last case of $\bar{M}$), we clearly cannot choose the same strand for $t$ and $u$, except if $|R| = 3$. Furthermore, we can clearly only minimize over new inner strands if such strands are still available. If $|R| \leq 3$, there may only be one available strand, or none at all, in which case the energy contribution is 0. We omit these edge cases in the presentation of the algorithm to maintain readability.

## B.2  Time complexity for Section 4.2

We need table entries for each possible configuration of remaining number of occurrences and for specifying the remaining number of bases on the leftmost and rightmost strand. Using $n := \max_{r \in R} |r|$, we bound the number of table entries by

$$n^2 \cdot \max_{s_1, \ldots, s_p : s_1 + \ldots + s_p = m} \prod_{i=1}^{p} s_p \leq n^2 \cdot \left(\frac{m}{p}\right)^p$$

The running time for computing one table entry is dominated, as for the previous section, by the last case. We need to iterate over $\mathcal{O}((\frac{m}{p})^p)$ configurations to split our region into two strand sets, $p$ lengths to determine the length of the strand on which we split and $n$ positions for the index of the split. We finally obtain a running time of $\mathcal{O}((\frac{m}{p})^{2p} \cdot n^3 \cdot p)$, which is an XP algorithm parametrized by $p$.

## B.3    Proof for the connectivity in Section 4.4

Analogous to Section 4.1, we define an interval $[m, s_i, r_j, c]$ and a relation $[m', t_k, u_\ell, c'] \preccurlyeq [m, s_i, r_j, c]$ if and only if $m' < m - 1$ or $m' = m - 1 \wedge (s = t \vee r = u)$ or $m' = m \wedge s = t \wedge r = u \wedge i \le k \wedge \ell \le j$. Since we just change the representation of our set $R$ to an integer $m$, the correctness of the algorithm can be shown by the same arguments as for the exponential algorithm. We only show here that the connectivity specifier $c \in \{1, 2\}$ actually enforces connectivity. For this, we introduce the following notation: $\gamma(m, s_i, r_j)$ means that the MFE structure computed by $M_{m, s_i, r_j, 1}$ is connected, and $\bar{\gamma}(m, s_i, r_j)$ means that the MFE structure computed by $M_{m, s_i, r_j, 2}$ is either connected or consists of two connected components, one containing $s$ and one containing $r$. In other words, adding a base pair between $s$ and $r$ to such a structure will make it connected. Let $[m] := \{1, ..., m\}$.

▶ **Lemma 28.** $\gamma(m', s_i, r_j) \wedge \bar{\gamma}(m', s_i, r_j)$ *for all* $m' \in [m]$, $s, r \in R$, $i \le |s|$ *and* $j \le |r|$.

**Proof.** Clearly, a secondary structure over an interval with $m = 1$ is always connected, i.e. $\gamma(1, t_k, t_\ell)$ and $\bar{\gamma}(1, t_k, t_\ell)$ hold for any valid $t, k, \ell$. By induction over $\preccurlyeq$, assume that $\gamma(m', t_k, u_\ell)$ and $\bar{\gamma}(m', t_k, u_\ell)$ for any $[m', t_k, u_\ell, c'] \preccurlyeq [m, s_i, r_j, c]$. We show $\gamma(m, s_i, r_j)$ and $\bar{\gamma}(m, s_i, r_j)$. By case distinction:

- **Case 1:** $s_i \notin S$. If $i + 1 \le |s|$, the structure is connected by assumption. Else, if $c = 2$, we need that a connection between $s$ and $r$ would make the structure connected. Indeed, by assumption, $[m - 1, t_1, r_j, 1]$ is connected, and together with a base pair between $s$ and $r$, all strands are in one connected component. If $c = 1$, $s$ and $r$ are not yet connected and we do not connect them with the last possible base $s_{|s|}$, thus no connected secondary structure with these constraints exists.
- **Case 2:** $\{s_i, r_j\} \in S$. By hypothesis, the structure for $]m, s_i, r_j, 2[$ would be connected together with a base pair between $s$ and $r$, thus the structure for $[m, s_i, r_j, c]$ is connected.
- **Case 3:** $\{s_i, t_k\} \in S$ for some $t_k$ in the region. By assumption and the base pair $\{s_i, t_k\}$, the strands from $s$ to $t$ are connected. If $c = 1$, then by assumption $]m - m' + 1, t_k, r_{j+1}, 1[$ is connected and thus all the structure is connected. For $c = 2$, assume a connection between $s$ and $r$. Now by the fact that $s$ is connected to $t$ and transitivity, $r$ is connected to $t$. We can apply our induction hypothesis to conclude that the substructure for the strands from $t$ to $r$ is connected, and by that, the complete structure is connected. ◀

We now argue (somewhat informally) why there cannot be a better connected secondary structure that the algorithm ignores. Assume that the last case of the $\bar{M}$ equation is defined as for 4.1, that is, we minimize over the two next inner strands. Any structure that uses this case cannot be connected (as the component including $s$ and $r$ has no way of being connected to the component including the inner strands).

Assume also that when minimizing over strands, we lift the connectivity requirement ($c = 1$). In any secondary structure than can be obtained by at some point (at interval $[m, s_i, r_j, 2]$) minimizing over a strand with $c = 2$ but not with $c = 1$, we know that the chosen inner strand (say $t$) is not connected to $r$ in the constructed secondary structure restricted to the region from $s_i$ to $r_j$. Since the outer region before $s_i$ and after $r_j$ does not contain any base of strand $t$, strand $t$ will not be connected to $r$ in the complete structure.

So, after applying these changes to the DP, we cannot achieve a better connected secondary structure than before. The DP is now almost equivalent to the DP in Section 4.1, with representing the set $R$ by a natural number $m$. We can thus repeat the correctness proof of section Section 4.1 to show that any (connected) secondary structure is covered by the equations, and thus the output of our DP is optimal.

## C    Appendix for Section 5

Figure 10 shows the probability of interior, exterior-homogeneous and exterior-heterogeneous base pairs for all pairs of TR, from $m = 2$ to 4. We can observe that the probabilities vary a lot and highly depend on the interacting triplets. Usually, internal and exterior-homogeneous base pairs behave similarly. One can also see that the probability of heterogeneous base pairs slightly increases with increasing $m$. On the other hand, the probability of observing interior base pairs is slightly decreasing.



Two interacting strands.

Three interacting strands.

Four interacting strands.

**Figure 10** Interaction profiles for pairs of triplets in the heterogeneous soup model.