

# Tree Diet: Reducing the Treewidth to Unlock FPT Algorithms in RNA Bioinformatics

Bertrand Marchand  

LIX CNRS UMR 7161, Ecole Polytechnique, Institut Polytechnique de Paris, Palaiseau, France  
LIGM, CNRS, Univ Gustave Eiffel, F77454 Marne-la-vallée, France

Yann Ponty<sup>1</sup>  

LIX CNRS UMR 7161, Ecole Polytechnique, Institut Polytechnique de Paris, Palaiseau, France

Laurent Bulteau<sup>1</sup>  

LIGM, CNRS, Univ Gustave Eiffel, F77454 Marne-la-vallée, France

---

## Abstract

Hard graph problems are ubiquitous in Bioinformatics, inspiring the design of specialized Fixed-Parameter Tractable algorithms, many of which rely on a combination of tree-decomposition and dynamic programming. The time/space complexities of such approaches hinge critically on low values for the treewidth  $tw$  of the input graph. In order to extend their scope of applicability, we introduce the TREE-DIET problem, *i.e.* the removal of a minimal set of edges such that a given tree-decomposition can be slimmed down to a prescribed treewidth  $tw'$ . Our rationale is that the time gained thanks to a smaller treewidth in a parameterized algorithm compensates the extra post-processing needed to take deleted edges into account.

Our core result is an FPT dynamic programming algorithm for TREE-DIET, using  $2^{O(tw)}n$  time and space. We complement this result with parameterized complexity lower-bounds for stronger variants (e.g., NP-hardness when  $tw'$  or  $tw - tw'$  is constant). We propose a prototype implementation for our approach which we apply on difficult instances of selected RNA-based problems: RNA design, sequence-structure alignment, and search of pseudoknotted RNAs in genomes, revealing very encouraging results. This work paves the way for a wider adoption of tree-decomposition-based algorithms in Bioinformatics.

**2012 ACM Subject Classification** Theory of computation  $\rightarrow$  Dynamic programming; Theory of computation  $\rightarrow$  Parameterized complexity and exact algorithms; Applied computing  $\rightarrow$  Bioinformatics

**Keywords and phrases** RNA, treewidth, FPT algorithms, RNA design, structure-sequence alignment

**Digital Object Identifier** 10.4230/LIPIcs.WABI.2021.7

**Related Version** *Full Version:* <https://hal.inria.fr/hal-03206132/>

**Supplementary Material** *Software (Source Code):* <https://gitlab.inria.fr/amibio/tree-diet>  
archived at `swh:1:dir:14ddd3079b8dca1f2cbea271a110cfb5592d23c`

**Funding** This work was supported by the French Agence Nationale de la Recherche through the Decrypted collaborative project (ANR-19-CE30-0021).

**Acknowledgements** The authors would like to thank Julien Baste for pointing out prior work on treewidth modulators, and providing valuable input regarding vertex deletion problems.

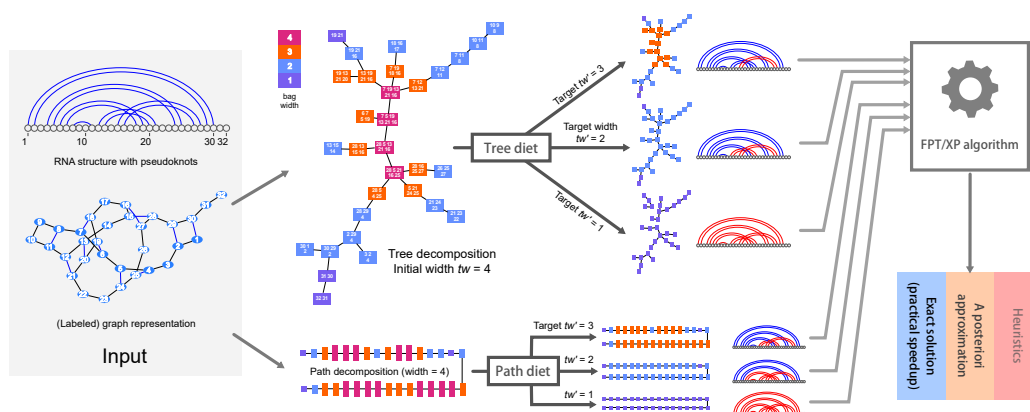
## 1 Introduction

Graph models and parameterized algorithms are found at the core of a sizable proportion of algorithmic methods in bioinformatics addressing a wide array of subfields, spanning sequence processing [48], structural bioinformatics [50], comparative genomics [9], phylogenetics [2], and further examples that can be found in a review by Bulteau and Weller [10]. RNA bioinformatics is no exception, with the prevalence of the secondary structure, an outer

---

<sup>1</sup> To whom correspondence should be addressed.





■ **Figure 1** General description of our approach and rationale. Starting from a structured instance, e.g. an RNA structure with pseudoknots, our tree-diet/path-diet algorithms extract simplified tree/path decompositions, having prescribed target width  $tw'$ . Those can be used within existing parameterized algorithms to yield efficient heuristics, a *a posteriori* approximations or even exact solutions.

planar graph [47], as an abstraction of RNA conformations, and the notable utilization of graph models to represent complex topological motifs called pseudoknots [49], inducing the hardness of several tasks, such as structure prediction [1, 29, 37], structure alignment [5], or structure/sequence alignment [32]. Such motifs are functionally important and conserved, as witnessed by their presence in the consensus structure of 336 RNA families in the 14.5 edition of the RFAM database [23]. Moreover, methods in RNA bioinformatics [36] are increasingly considering non-canonical base pairs and modules [25, 31], further increasing the density of RNA structural graphs and outlining the need for scalable algorithms.

A parameterized complexity approach can be used to circumvent the frequent NP-hardness of relevant problems. It generally considers one or several parameters, whose values are naturally bounded (or much smaller than the input size) within real-life instances. Once relevant parameters have been identified, one aims to design a Fixed Parameter Tractable (FPT) algorithm, having polynomial complexity for any fixed value of the parameter, and reasonable dependency on the parameter value. The treewidth is a classic parameter for FPT algorithms, and intuitively captures a notion of distance of the input to a tree. It is popular in bioinformatics due to the existence of efficient heuristics [19, 8] for computing tree-decompositions of reasonable treewidth. Given a tree-decomposition, many combinatorial optimization tasks can be solved using dynamic programming (DP), in time/space complexities that remain polynomial for any fixed treewidth value. Resulting algorithms remain correct upon (almost) arbitrary modifications of the objective function parameters, and can be adapted to study statistical properties of search spaces through changes of algebra.

Unfortunately, the existence of a parameterized (or FPT) algorithm does not necessarily imply that of a practically-efficient implementation, even when the parameter takes low typical values. Indeed, the dependency of the complexity on the treewidth may be prohibitive, both in terms of time and memory requirements. This limitation is particularly obvious while searching and aligning structured RNAs, giving rise to an algorithmic problem called the RNA structure-sequence alignment [39, 21, 32], for which the best known exact algorithm is in  $\Theta(n \cdot m^{tw+1})$ , where  $n$  is the structure length,  $m$  the sequence/window, and  $tw$  is the treewidth of the structure (inc. backbone). Similar complexities hold for problems that can be expressed as (weighted) constraint satisfaction problems, with  $m$  representing the

cardinality of the variable domains. Such frameworks are frequently used for molecular design, both in proteins [44] and RNA [51], and may require the consideration of tree-widths up to 20 or more [20].

In this paper, we investigate a pragmatic strategy to increase the practicality of parameterized algorithms based on the treewidth parameter [6]. We put our instance graphs on a diet, *i.e.* we introduce a preprocessing that reduces their treewidth to a prescribed value by removing a minimal cardinality set of edges. As discussed previously, the practical complexity of many algorithms greatly benefits from the consideration of simplified instances, having lower treewidth. Moreover, specific countermeasures for errors introduced by the simplification can sometimes be used to preserve the correctness of the algorithm. For instance, searching structured RNAs using RNA structure-sequence alignment [39], an iterated filtering strategy could use instances of increasing treewidth to restrict potential hits, weeding them early so that a – costly – full structure is reserved to (quasi-)hits. This strategy could remain exact while saving substantial time. Alternative countermeasures could be envisioned for other problems, such as a rejection approach to correct a bias introduced by simplified instances in RNA design.

After stating our problem(s) in Section 2, we study in Section 3 the parameterized complexity of the GRAPH-DIET problem, the removal of edges to reach a prescribed treewidth. We propose, in Section 4, a practical Dynamic Programming FPT algorithm for TREE-DIET, along with possible further optimizations for PATH-DIET, two natural simplifications of the GRAPH-DIET problem, where a tree (resp. path) decomposition is provided as input and used as a guide. Finally, we show in Section 5 how our algorithm can be used to extract hierarchies of graphs/structural models of increasing complexity to provide alternative sampling strategies for RNA design, and speed-up the search for pseudoknotted non-coding RNAs. We conclude in Section 6 with future considerations and open problems.

## 2 Statement of the problem(s) and results

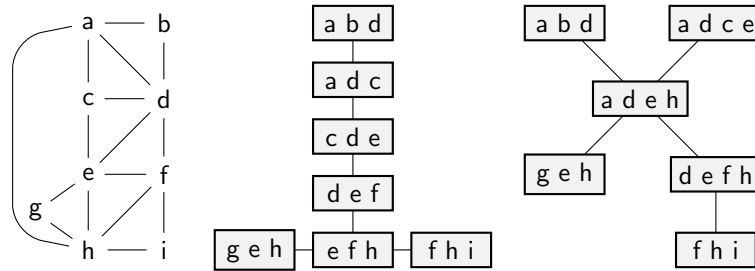
A *tree-decomposition*  $\mathcal{T}$  (over a set  $V$  of vertices) is a tree whose nodes are subsets of  $V$ , known as bags. The bags containing any  $v \in V$  induce a (connected) subtree of  $\mathcal{T}$ . A *path-decomposition* is a tree-decomposition whose underlying tree  $\mathcal{T}$  is a path. The *width* of  $\mathcal{T}$  (denoted  $w(\mathcal{T})$ ) is the size of its largest bag minus 1. An edge  $\{u, v\}$  is *visible* in  $\mathcal{T}$  if some bag contains both  $u$  and  $v$ , otherwise it is *lost*.  $\mathcal{T}$  is a *tree-decomposition of  $G$*  if all edges of  $G$  are visible in  $\mathcal{T}$ . The *treewidth* of a graph  $G$  is the minimum width over all tree-decompositions of  $G$ .

► **Problem (GRAPH-DIET).** *Given a graph  $G = (V, E)$  of treewidth  $tw$ , and an integer  $tw' < tw$ , find a tree-decomposition over  $V$  of width at most  $tw'$  losing a minimum number of edges from  $G$ .*

A *tree-diet* of  $\mathcal{T}$  is any tree-decomposition  $\mathcal{T}'$  obtained by removing vertices from the bags of  $\mathcal{T}$ .  $\mathcal{T}'$  is a  *$d$ -tree-diet* if  $w(\mathcal{T}') \leq w(\mathcal{T}) - d$ .

► **Problem (TREE-DIET).** *Given a graph  $G$ , a tree-decomposition  $\mathcal{T}$  of  $G$  of width  $tw$ , and an integer  $tw' < tw$ , find a  $(tw - tw')$ -tree-diet of  $\mathcal{T}$  losing a minimum number of edges.*

Note that for TREE-DIET,  $\mathcal{T}$  does not have to be optimal, so the width  $tw$  of the input tree decomposition might be larger than the actual treewidth of  $G$ , thus TREE-DIET can be used to reduce the width of *any* input decomposition. We define BINARY-TREE-DIET and PATH-DIET analogously, where  $\mathcal{T}$  is restricted to be a binary tree (respectively, a path). An example of an instance of GRAPH-DIET and of TREE-DIET are given in Figure 2.



■ **Figure 2** Illustrations for the GRAPH-DIET and TREE-DIET problems. Given a graph  $G$  on the left (treewidth 3), an optimal solution for GRAPH-DIET, with target treewidth 2, yields the tree-decomposition in the middle (edge  $ah$  is lost). On the other hand, any 1-tree-diet for the tree-decomposition on the right loses at least 3 edges.

## Parameterized Complexity in a Nutshell

The basics of parameterized complexity can be loosely defined as follows (see [17] for the formal background). A *parameter*  $k$  for a problem is an integer associated with each instance which is expected to remain small in practical instances (especially when compared to the input size  $n$ ). An exact algorithm, or the problem it solves, is FPT if it takes time  $f(k)\text{poly}(n)$ , and XP if it takes time  $n^{g(k)}$  (for some functions  $f, g$ ). Under commonly accepted conjectures (see for instance [15] for details), W[1]-hard problems may not be FPT, and Para-NP-hard problems (NP-hard even for some fixed value of  $k$ ) are not FPT nor XP.

### 2.1 Our results

Our results are summarized in Table 1. Although the GRAPH-DIET problem would give the most interesting tree-decompositions in theory, it seems unlikely to admit efficient algorithms in practice (see Section 3).

Thus we focus on the TREE-DIET relaxation, where an input tree-decomposition is given, which we use as a guide/restriction towards a thinner tree-decomposition. Seen as an additional constraint, it makes the problem harder (the case  $tw' = 1$  becomes NP-hard, Theorem 3, although for GRAPH-DIET it corresponds to the SPANNING TREE problem and is polynomial). With parameter  $tw$  however, it does help reduce the search space. In Theorem 10 we give an  $O((6\Delta)^{tw}\Delta^2n)$  Dynamic Programming algorithm, where  $\Delta$  is the maximum number of children of any bag in the tree-decomposition. This algorithm can thus be seen as XP in general, but FPT on bounded-degree tree-decompositions (e.g. binary trees and paths). This is not a strong restriction, since the input tree may safely and efficiently be transformed into a binary one (see Appendix A for more details). Moreover, the duplications of bags which are used in the conversion may only decrease the number of lost edges incurred by TREE-DIET.

We also consider the case where the treewidth needs to be reduced by  $d = 1$  only, this without constraining the source treewidth. We give a polynomial-time algorithm for PATH-DIET in this setting (Theorem 13) which generalizes into an XP algorithm for larger values of  $d$ , noting that an FPT algorithm for  $d$  is out of reach by Theorem 5. We also show that the problem is Para-NP-hard if the tree degree is unbounded (Theorem 4).

■ **Table 1** Parameterized results for our problems. Algorithm complexities are given up to polynomial time factors ( $O^*$  notation),  $\Delta$  denotes the maximum number of children in the input tree-decomposition. (\*) see Theorem 2 statement for a more precise formulation.

Parameter Problem	Source treewidth $tw$		Target treewidth $tw'$	Difference $d = tw - tw'$	
GRAPH-DIET	<i>open</i>		Para-NP-hard $tw' = 2$ EDP( $K_4$ ) [18]	Para-NP-hard* $d = 1$ Theorem 2	
TREE-DIET	XP $O^*((6\Delta)^{tw})$ Theorem 10	FPT <i>open</i>	Para-NP-hard $tw' = 1$ Theorem 3	Para-NP-hard $d = 1$ Theorem 4	
BINARY-TREE-DIET	FPT $O^*(12^{tw})$ Theorem 10			W[1]-hard Theorem 5	XP <i>open</i>
PATH-DIET					XP $O^*(tw^d)$ Theorem 13

### 3 Algorithmic Limits: Parameterized Complexity Considerations

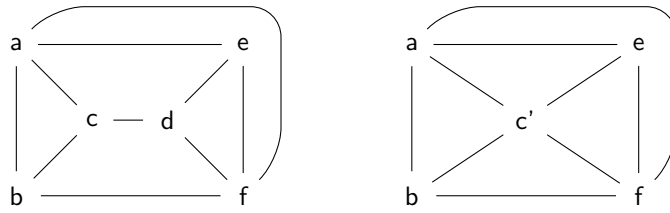
GRAPH-DIET can be seen as a special case of the EDGE DELETION PROBLEM (EDP) for the family of graphs  $\mathcal{H}$  of treewidth  $tw'$  or less: given a graph  $G$ , remove as few edges as possible to obtain a graph in  $\mathcal{H}$ . Such edge modification problems are more often parameterized by the number  $k$  of edited edges (see [14] for a complete survey). Given our focus on increasing the practicality of treewidth-based algorithms in bioinformatics, we restrict our focus to treewidth related parameters  $tw$ ,  $tw'$  and  $d = tw - tw'$ .

Considering the target treewidth  $tw'$ , we note that EDP is NP-hard when  $\mathcal{H}$  is the family of treewidth-2 graphs [18], namely  $K_4$ -free graphs, hence the notation EDP( $K_4$ ). It follows that GRAPH-DIET is Para-NP-hard for the target treewidth parameter  $tw'$ .

#### 3.1 Graph-Diet: practical solutions seem unlikely

For a combination of the parameters  $tw'$  and  $k$ , we could imagine graph minor theorems yielding parameterized algorithms “for free”, as it is often the case with treewidth-based problems. In this respect, GRAPH-DIET corresponds to deciding if a graph  $G$  belongs to the family of graphs having treewidth  $tw'$ , augmented by  $k$  additional edges, denoted as  $\text{Treewidth-}tw'+ke$  since its introduction by Cai [12]. If this family were minor-closed, polynomial minor-free-testing [27, 34] would yield an FPT algorithm. However, this is not the case: for some graphs in the family, an edge contraction yields a graph  $G'$  not in  $\text{Treewidth-}tw'+ke$ , as illustrated by Figure 3.

Regarding the source graph treewidth  $tw$ , a theoretical approach could be via Courcelle’s Theorem [13] and Monadic Second Order (MSO) formulas: it suffices to express the property of being in  $\text{Treewidth-}tw'+ke$  using MSO to prove that GRAPH-DIET is FPT for  $tw$ . We presume this is feasible (the main brick being minor testing, which is expressible in MSO), however it is not clear whether this is doable with formulas independent of  $k$ , in order to obtain an algorithm for the treewidth alone. In any case, this approach would probably not yield practical algorithms. Indeed, Courcelle’s theorem typically lead to running times involving towers of exponentials on the relevant parameters, so we do not investigate it further within the scope of this paper.



■ **Figure 3** A graph  $G$  (left) with treewidth 3. Deleting edge  $cd$  gives treewidth 2, implying that  $G \in \text{Treewidth}2 + 1e$ . However, if one contracts edge  $cd$ , then the resulting graph (right) has treewidth 3, and deleting any single edge does not decrease the treewidth. This example shows that the graph family  $\text{Treewidth } 2+1e$  is not minor-closed.

Another meta-theorem by Cai [11] may yield an FPT algorithm if  $\text{Treewidth-}tw'+ke$  can be described through a finite number forbidden induced subgraphs, but again  $k$  would most likely be a parameter as well. On a related note, it is worth noting that Edge Deletion to other graph classes (interval, permutation, ...) does admit efficient algorithms when parameterized by the treewidth alone [35], painting a contrasted picture.

The vertex deletion equivalent of GRAPH-DIET, where one asks for a minimum subset of vertices to remove to obtain a given treewidth, is known as a TREEWIDTH MODULATOR. This problem has been better-studied than its edge-deletion counterpart [16], and has been shown to be FPT for the treewidth [3], with a reasonable dependency in the parameter. However, it is currently unclear how this can be adapted into an edge deletion algorithm.

Overall an FPT algorithm for GRAPH-DIET does not seem out of reach, and could result from one of the above-mentioned meta-theorems. However it seems unlikely to induce a “practical” exact algorithms. Indeed, any algorithm for GRAPH-DIET can be used to compute the TREEWIDTH of an arbitrary graph, for which current state-of-the-art exact algorithms require time in  $tw^{O(tw^3)}$  [6]. We thus have the following conjecture, which motivates the TREE-DIET relaxation of the problem.

► **Conjecture 1.** GRAPH-DIET is FPT for the source treewidth parameter ( $tw$ ), but no algorithm with single-exponential running time exists.

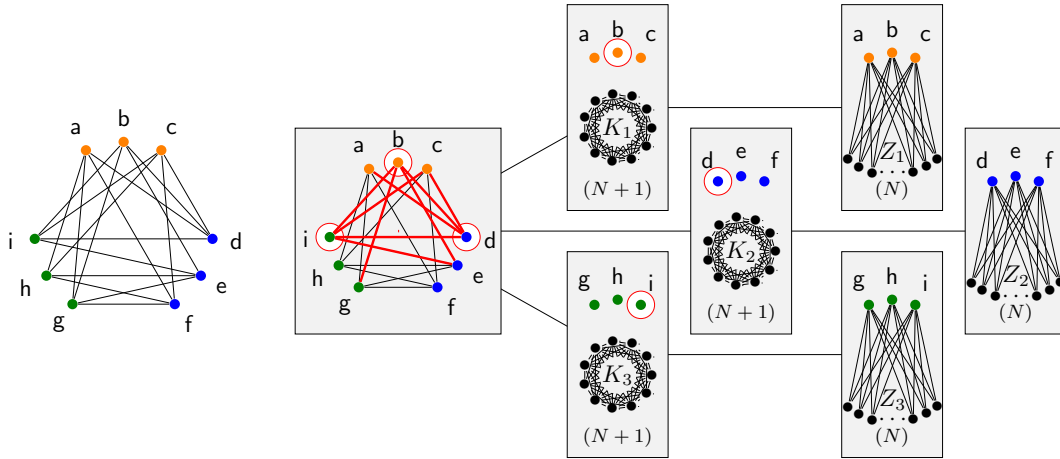
Finally, for parameter  $d$ , any polynomial-time algorithm for constant  $d$  would allow to compute the treewidth of any graph in polynomial time. Since treewidth is NP-hard we have the following result

► **Theorem 2.** There is no XP algorithm for GRAPH-DIET with parameter  $d$  unless  $P=NP$ .

### 3.2 Lower Bounds for Tree-Diet

Parameters  $tw'$  and  $d$  would be the most interesting in practice, since parameterized algorithms would be efficient for small diets or small target treewidth. However, we prove strong lower-bounds for TREE-DIET on each of these parameters, leaving very little hope for parameterized algorithms (we thus narrow down the possible algorithms to the combined parameter  $tw'+d$ , i.e.  $tw$ , see Section 4). Only XP for parameter  $d$  when  $\mathcal{T}$  has a constant degree remains open (cf. Table 1).

► **Theorem 3.** TREE-DIET and PATH-DIET are Para-NP-hard for the target treewidth parameter  $tw'$  (NP-hard for  $tw' = 1$ ).



■ **Figure 4** Reduction for Theorem 4 showing that TREE-DIET is NP-hard even for  $d = 1$ , from a graph  $G$  (left) with  $k = 3$  and  $n = 3$  to a graph  $G'$  (right, given by its tree-decomposition of width  $N + n + 1$ ): a 1-tree-diet for  $G'$  amounts to selecting a  $k$ -clique in the root bag, i.e. in  $G$ .

**Proof.** By reduction from the NP-hard problem SPANNING CATERPILLAR TREE [42]: given a graph  $G$ , does  $G$  has a spanning tree  $C$  that is a caterpillar? Given  $G = (V, E)$  with  $n = |V|$ , we build a tree-decomposition  $\mathcal{T}$  of  $G$  consisting of  $n - 1$  bags containing all vertices (the width of the decomposition is therefore  $n - 1$ ) connected in a path. Then  $(G, \mathcal{T})$  admits a tree-diet to treewidth 1 with  $n - 1$  visible edges if, and only if,  $G$  admits a caterpillar spanning tree. Indeed, the subgraph of  $G$  with visible edges must be a graph with pathwidth 1, i.e. a caterpillar [30]. With  $n - 1$  visible edges, the caterpillar connects all  $n$  vertices together, i.e. it is a spanning tree. ◀

► **Theorem 4.** TREE-DIET is Para-NP-hard for parameter  $d$ . More precisely, it is W[1]-hard for parameter  $\Delta$ , the degree of  $\mathcal{T}$ , even when  $d = 1$ .

**Proof.** By reduction from MULTI-COLORED CLIQUE. Consider a  $k$ -partite graph  $G = (V, E)$  with  $V = \bigcup_{i=1}^k V_i$ . We assume that  $G$  is regular (each vertex has degree  $\delta$  and that each  $V_i$  has the same size  $n$  (MULTI COLORED CLIQUE is W[1]-hard under these restrictions [17, 15]). Let  $L := \delta k - \binom{k}{2}$  and  $N = \max\{|V|, L + 1\}$ . We now build a graph  $G'$  and a tree-decomposition  $\mathcal{T}'$ : start with  $G' := G$ . Add  $k$  independent cliques  $K_1, \dots, K_k$  of size  $N + 1$ . Add  $k$  sets of  $N$  vertices  $Z_i$  ( $i \in [k]$ ) and, for each  $i \in [k]$ , add edges between each  $v \in V_i$  and each  $z \in Z_i$ . Build  $\mathcal{T}$  using  $2k + 1$  bags  $T_0, T_{1,i}, T_{2,i}$  for  $i \in [k]$ , such that  $T_0 = V$ ,  $T_{1,i} = V_i \cup K_i$  and  $T_{2,i} = V_i \cup Z_i$ . The tree-decomposition is completed by connecting  $T_{2,i}$  to  $T_{1,i}$  and  $T_{1,i}$  to  $T_0$  for each  $i \in [k]$ . Thus,  $\mathcal{T}$  is a tree-decomposition of  $G'$  with  $\Delta = k$  and maximum bag size  $n + N + 1$  (vertices of  $V$  induce a size-3 path in  $\mathcal{T}$ , other vertices appear in a single bag, edges of  $G$  appear in  $T_0$ , edges of  $K_i$  in  $T_{1,i}$ , and finally edges between  $V_i$  and  $Z_i$  appear in  $T_{2,i}$ ). The following claim completes the reduction:

$\mathcal{T}$  has a 1-tree-diet losing at most  $L$  edges from  $G' \Leftrightarrow G$  has a  $k$ -clique.

◀ Assume  $G$  has a  $k$ -clique  $X = \{x_1, \dots, x_k\}$  (with  $x_i \in V_i$ ). Build  $\mathcal{T}'$  by removing each  $x_i$  from bags  $T_0$  and  $T_{1,i}$ . Then  $\mathcal{T}'$  is a 1-tree-diet of  $\mathcal{T}$ . There are no edges lost by removing  $x_i$  from  $T_{1,i}$  (since  $x_i$  is not connected to  $K_i$ ), and the edges lost in  $T_0$  are all edges of  $G$  adjacent to any  $x_i$ . Since  $X$  forms a clique and each  $x_i$  has degree  $\delta$ , there are  $L = k\delta - \binom{k}{2}$  such edges.

$\Rightarrow$  Consider a 1-tree-diet  $\mathcal{T}'$  of  $\mathcal{T}$  losing  $L$  edges. Since each bag  $T_{1,i}$  has maximum size,  $\mathcal{T}'$  must remove at least one vertex  $x_i$  in each  $T_{1,i}$ . Note that  $x_i \in V_i$  (since removing  $x_i \in K_i$  would lose at least  $N \geq L + 1$  edges). Furthermore,  $x_i$  may not be removed from  $T_{2,i}$  (otherwise  $N$  edges between  $x_i$  and  $Z_i$  would be lost), so  $x_i$  must also be removed from  $T_0$ . Let  $K$  be the number of edges in  $G[\{x_1 \dots x_k\}]$ . The total number of lost edges in  $T_0$  is  $\delta k - K$ . Thus, we have  $\delta k - K \leq L$  and  $K \geq \binom{k}{2}$ :  $\{x_1, \dots, x_k\}$  form a  $k$ -clique of  $G$ .  $\blacktriangleleft$

► **Theorem 5.** PATH-DIET is  $W[1]$ -hard for parameter  $d$ .

## 4 FPT Algorithm

### 4.1 For general tree-decompositions

We describe here a  $O(3^{tw}n)$ -space,  $O(\Delta^{tw+2} \cdot 6^{tw}n)$ -time dynamic programming algorithm for the TREE-DIET problem, with  $\Delta$  and  $tw$  being respectively the maximum number of children of a bag in the input tree-decomposition and its width. On *binary* tree-decompositions (where each bag has at most 2 children), it yields a  $O(3^{tw}n)$ -space  $O(12^{tw}n)$ -time FPT algorithm.

#### 4.1.1 Coloring formulation

We aim at solving the following problem: given a tree-decomposition  $\mathcal{T}$  of width  $tw$  of a graph  $G$ , we want to remove vertices from the bags of  $\mathcal{T}$  to reach a target width  $tw'$  while *losing* as few edges from  $G$  as possible. We tackle the problem through an equivalent *coloring* formulation: our algorithm will assign a color to each occurrence of a vertex in the tree decomposition. We work with three colors: red (r), orange (o), and green (g). Green means that the vertex is kept in the bag, while orange and red means removal of the vertex. An edge is thus visible within a bag when both its ends are green. It is lost if there is no bag where it is visible. To ensure equivalence with the original problem, the colors will be assigned following local rules, which we now describe.

We first root the tree-decomposition arbitrarily.

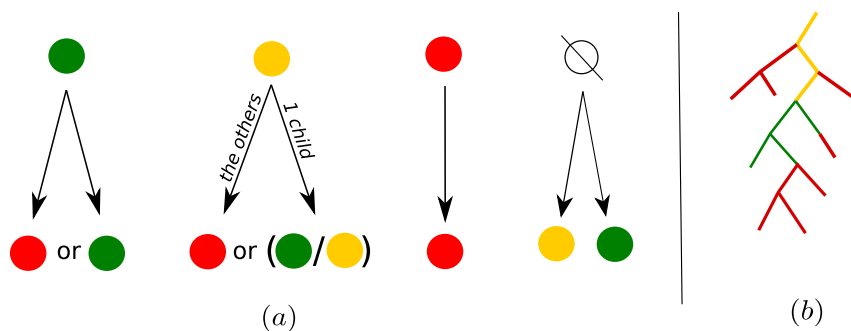
► **Definition 6.** A coloring of vertices in the bags of the decomposition is said to be valid if it follows the following rules:

- A vertex of a bag not present in its parent may be green or orange (R1)
- A green vertex in a bag may be either green or red in its children (R2)
- A red vertex in a bag must stay red in its children (R3)
- An orange vertex in a bag has to be either orange or green in exactly one child (unless there is no child with this vertex), and must be red in the other children (R4)

These rules are summarized in Figure 5 (a).

When going down the tree, a green vertex may only stay green or permanently become red. An immediate consequence of these rules is therefore that the green occurrences of a given vertex form a (possibly empty) connected subtree. Informally, orange vertices are locally absent but “may potentially be found further down the tree”, while red vertices are removed from both the current bag and its entire subtree. Figure 5 (b) shows an example sketch for a valid coloring of the occurrences of a given vertex in the tree-decomposition. A vertex may only be orange along a path starting from its highest occurrence in the tree, with any part branching off that path entirely red. It ends at the top of a (potentially empty) green subtree, whose vertices may also be parents to entirely red subtrees.





■ **Figure 5** (a) Color assignment rules for vertices, when going down-tree. (b) Sketch of the general pattern our color assignment rules create on  $\mathcal{T}_u$ , the subtree of bags containing a given vertex  $u$ .

We will now more formally prove the equivalence of the coloring formulation to the original problem. Let us first introduce two definitions. Given a valid coloring  $\mathcal{C}$  of a tree-decomposition of  $G$ , an edge  $(u, v)$  of  $G$  is said to be *realizable* if there exists a bag in which both  $u$  and  $v$  are green per  $\mathcal{C}$ . Given an integer  $d$ , a coloring  $\mathcal{C}$  of  $\mathcal{T}$  is said to be  *$d$ -diet-valid* if removing red/orange vertices reduces the width of  $\mathcal{T}$  from  $w(\mathcal{T})$  to  $w(\mathcal{T}) - d$ .

► **Proposition 7.** *Given a graph  $G$ , a tree-decomposition  $\mathcal{T}$  of width  $tw$ , and a target width  $tw' < tw$ , The TREE-DIET problem is equivalent to finding a  $(tw - tw')$ -diet-valid coloring  $\mathcal{C}$  of  $\mathcal{T}$  allowing for a number of realizable edges in  $G$  as large as possible.*

**Proof.** Given a  $(tw - tw')$ -tree-diet of  $\mathcal{T}$ , specifying which vertices are removed from which bags, we obtain a valid coloring  $\mathcal{C}$  for  $\mathcal{T}$  incurring the same number of lost (unrealizable) edges. To start with, a vertex  $u$  is colored green in the bags where it is not removed. By the validity of  $\mathcal{T}'$  as a decomposition, this set of bags forms a connected subtree, that we denote  $\mathcal{T}_u^g$ . We also write  $\mathcal{T}_u$  for the subtree of bags containing  $u$  in the original decomposition  $\mathcal{T}$ . If  $\mathcal{T}_u^g$  and  $\mathcal{T}_u$  do not have the same root, then  $u$  is colored orange on the the path in  $\mathcal{T}$  from the root of  $\mathcal{T}_u$  (included) and the root of  $\mathcal{T}_u^g$  (excluded). Vertex  $u$  is colored red in any other bag of  $\mathcal{T}_u$  not covered by these two cases. The resulting coloring follows rules (R1-4) and induces the same set of lost/non-realizable edges as the original  $(tw - tw')$ -tree-diet. Conversely, an equivalent  $(tw - tw')$ -tree-diet is obtained from a  $(tw - tw')$ -diet-valid coloring by removing red/orange vertices and keeping green ones. If a given vertex has no green occurrences, it is entirely removed from the tree decomposition and all its edges are *lost* (it becomes an isolated vertex). We may add it back to the tree decomposition by introducing a new bag containing only this vertex, which we connect arbitrarily to the tree decomposition. ◀

#### 4.1.2 Decomposition of the search space and sub-problems

Based on this coloring formulation, we now describe a dynamic programming scheme for the TREE-DIET problem. We work with sub-problems indexed by tuples  $(X_i, f)$ , with  $X_i$  a bag of the input tree decomposition and  $f$  a coloring of the vertices of  $X_i$  in green, orange or red (in particular,  $f^{-1}(g)$  denotes the green vertices of  $X_i$ , and similarly for  $o$  and  $r$ ).

Let us introduce some notations before giving the definition of our dynamic programming table. Given an edge  $(u, v)$  of  $G$ , realizable when coloring a tree-decomposition  $\mathcal{T}$  of  $G$  with  $\mathcal{C}$ , we write  $\mathcal{T}_{uv}^g$  the subtree of  $\mathcal{T}$  in which both  $u$  and  $v$  are green. We denote by  $\mathcal{T}_i$  the subtree of the decomposition rooted at  $X_i$ , and  $\mathcal{C}(i, f)$  the  $d$ -diet-valid colorings of  $\mathcal{T}_i$  agreeing with  $f$  on  $i$ , with  $d = tw - tw'$ . Our dynamic programming table is then defined as:

$$c(X_i, f) = \begin{cases} \max_{\mathcal{C} \in \mathcal{C}^{(i,f)}} \left| \left\{ \begin{array}{l} \text{Edges } (u, v) \text{ of } G, \text{ realizable within } \mathcal{T}_i \text{ colored with } \mathcal{C} \\ \text{such that } \mathcal{T}_{uv}^g \text{ is entirely contained strictly below } X_i \end{array} \right\} \right| & \text{if } f \text{ assigns green to at most } tw' + 1 \text{ vertices} \\ -\infty & \text{otherwise} \end{cases}$$

The cell  $c(X_i, f)$  therefore aggregates all edges realizable *strictly below*  $X_i$ . As we shall see through the recurrence relation below and its proof, edges with both ends green in  $X_i$  will be accounted for *above*  $X_i$  in  $\mathcal{T}$ .

We assume w.l.o.g that the tree-decomposition is rooted at an empty bag  $R$ . Given the definition of the table, the maximum number of realizable edges, compatible with a tree-diet of  $(tw - tw')$  to  $\mathcal{T}$ , can be found in  $c(R, \emptyset)$ .

The following theorem presents a recurrence relation obeyed by  $c(X_i, f)$  :

► **Theorem 8.** *For a bag  $X_i$  of  $\mathcal{T}$ , with children  $Y_1, \dots, Y_\Delta$ , we have:*

$$c(X_i, f) = \max_{m: f^{-1}(o) \rightarrow [1.. \Delta]} \left[ \sum_{1 \leq j \leq \Delta} \left( \max_{f'_j \in \text{compatible}(Y_j, f, m)} c(Y_j, f'_j) + |\text{count}(f, f'_j)| \right) \right]$$

with

- $m$ : a map from the orange vertices in  $X_i$  to the children of  $X_i$ . It decides for each orange vertex  $u$ , which child, among those which contain  $u$ , will color  $u$  orange or green; If there are no orange vertices in  $X_i$ , only the trivial empty map is considered.
- $\text{compatible}(Y_j, f, m)$ : the set of colorings of  $Y_j$  compatible with  $f$  on  $X_i$  and  $m$ ;
- $\text{count}(f, f'_j)$ : set of edges of  $G$  involving two vertices of  $Y_j$  green by  $f'_j$ , but such that one of them is either not in  $X_i$  or not green by  $f$ .

Note that  $\text{compatible}(Y_j, f, m)$  may contain colorings  $f'_j$  that colour too many vertices in  $Y_j$  in green to reach target width  $tw'$ . In that case  $c(Y_j, f'_j) = -\infty$ .

Theorem 8 relies on the following separation lemma for realizable edges under a valid coloring of a tree-decomposition. Recall that we suppose w.l.o.g that the tree-decomposition is rooted at an empty bag.

► **Lemma 9.** *An edge  $(u, v)$  of  $G$ , realizable in  $\mathcal{T}$  under  $\mathcal{C}$ , is contained in exactly one set of the form  $\text{count}(C_{|P}, C_{|X})$  with  $X$  a bag of  $\mathcal{T}$  and  $P$  its parent,  $C_{|P}, C_{|X}$  the restrictions of  $\mathcal{C}$  to  $P$  and  $X$ , respectively, and “count” defined as above. In addition,  $X$  is the root of the subtree of  $\mathcal{T}$  in which both  $u$  and  $v$  are green.*

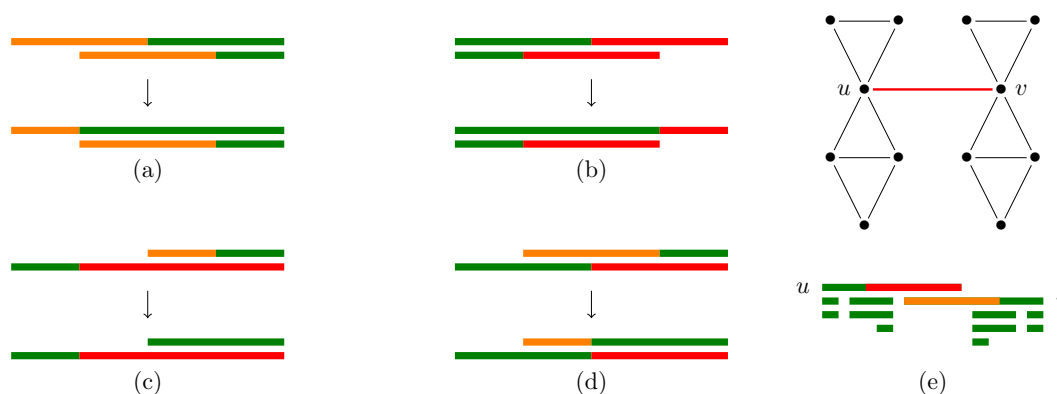
### Dynamic programming algorithm

The recurrence relation of Theorem 8 naturally yields a dynamic programming algorithm for the TREE-DIET problem, as stated below:

► **Theorem 10.** *There exists a  $O(\Delta^{tw+2} \cdot 6^{tw} \cdot n)$ -time,  $O(3^{tw} \cdot n)$ -space algorithm for the TREE-DIET problem, with  $\Delta$  the maximum number of children of a bag in the input tree-decomposition, and  $tw$  its width.*

► **Corollary 11.** BINARY-TREE-DIET ( $\Delta = 2$ ) admits an FPT algorithm for the  $tw$  parameter.

A pseudo-code implementation of the algorithm, using memoization, is included in Appendix B.



■ **Figure 6** Five cases where two vertices are deleted in the same bag with  $d = 1$ . Bags are points in the line, and an interval covering all bags containing  $v$  is drawn for each  $v$  (with an equivalent coloring, see Proposition 7). Cases (a) to (d) can be safely avoided by applying the given transformations. In the example for case (e), however, it is necessary to delete both vertices  $u$  and  $v$  from a central bag. It is sufficient to avoid cases (a) and (b) in order to obtain an XP algorithm for  $d$ .

## 4.2 For path decompositions

In the context of paths decompositions, we note that the number of removed vertices per bag can be limited to at most  $2d$  without losing the optimality. More precisely, we say that a coloring  $\mathcal{C}$  is  $d$ -simple if any bag has at most  $d$  orange and  $d$  red vertices. We obtain the following result, using transformations given in Figure 6.

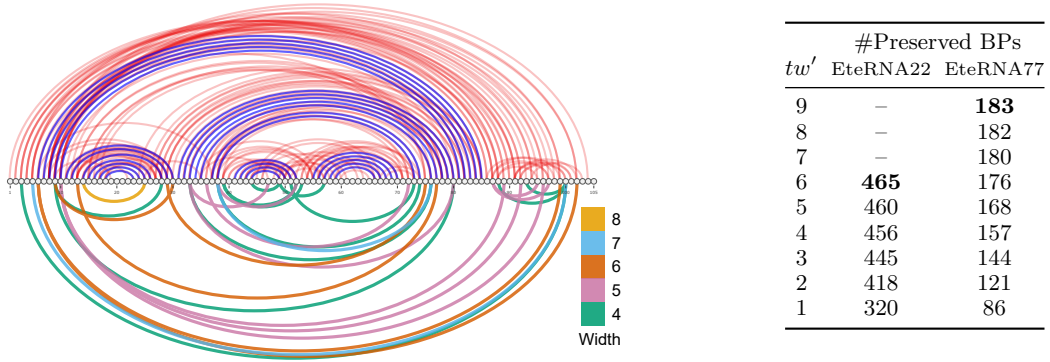
► **Proposition 12.** *Given a graph  $G$  and a path-decomposition  $\mathcal{T}$ , if  $\mathcal{C}$  is a  $d$ -diet-valid coloring of  $\mathcal{T}$  losing  $k$  edges, then  $\mathcal{T}$  has a  $d$ -diet-valid coloring that is  $d$ -simple, and loses at most  $k$  edges.*

Together with Proposition 7, this shows that it is sufficient to restrict our algorithm to  $d$ -simple colorings. (See also Figure 6). In particular, for any set  $X_i$ , choosing which  $\leq d$  vertices are orange and which  $\leq d$  are red, among the total of  $n$  vertices, is enough to fix a coloring. The number of such colorings is therefore bounded by  $O(tw^{2d})$ . Applying this remark to our algorithm presented in Section 4.1 yields the following result:

► **Theorem 13.** *PATH-DIET can be solved in  $O(tw^{2d}n)$ -space and  $O(tw^{4d}n)$ -time.*

## 5 Proofs of concept

We now illustrate the relevance of our approach, and the practicality of our algorithm for TREE-DIET, by using it in conjunction with FPT algorithms for three problems in RNA bioinformatics. We implemented in C++ the dynamic programming scheme described in Theorem 10 and Appendix B. Its main primitives are made available for Python scripting through pybind11 [22]. It actually allows to solve a generalized *weighted* version of TREE DIET, as explained in Appendix B. This feature allows to favour the conservation of important edges (e.g. RNA backbone) during simplification, by assigning them a much larger weight compared to other edges. Our implementation is freely available at <https://gitlab.inria.fr/amibio/tree-diet>.



■ **Figure 7** (Left) Target secondary structure (blue BPs), full set of disruptive base pairs (DPB; top) inferred by RNAPond on the Eterna77 puzzle, and subsets of DBPs (bottom) cumulatively removed by the tree-diet algorithm to reach prescribed treewidths. (Right) Number of BPs retained by our algorithm, targeting various treewidth values for the EteRNA22 and EteRNA77 puzzles.

## 5.1 Memory-parsimonious unbiased sampling of RNA designs

As a first use case for our simplification algorithm, we strive to ease the sampling phase of a recent method, called RNAPond [51], addressing RNA negative design. The method targets a set of base pairs  $S$ , representing a secondary structure of length  $n$ , and infers a set  $\mathcal{D}$  of  $m$  disruptive base pairs (DBPs) that must be avoided. It relies on a  $\Theta(k \cdot (n + m))$  time algorithm for sampling  $k$  random sequences (see Appendix C for details) after a preprocessing in  $\Theta(n \cdot m \cdot 4^{tw})$  time and  $\Theta(n \cdot 4^{tw})$  space. Here, the input consists of a graph  $G = ([1, n], S \cup \mathcal{D})$  and a tree decomposition  $\mathcal{T}$  of  $G$ , having width  $tw$ . In practice, the preprocessing largely dominates the overall runtime, even for large values of  $k$ , and its large memory consumption represents the main bottleneck.

This discrepancy in the complexities/runtimes of the preprocessing and sampling suggests an alternative strategy: relaxing the set of constraints to  $(S', \mathcal{D}')$ , with  $(S' \cup \mathcal{D}') \subset (S \cup \mathcal{D})$ , and compensating it through a rejection of sequences violating constraints in  $(S, \mathcal{D}) \setminus (S', \mathcal{D}')$ . The relaxed algorithm would remain unbiased, while the average-case time complexity of the rejection algorithm would be in  $\Theta(k \cdot \bar{q} \cdot (n + m))$  time, where  $\bar{q}$  represents the relative increase of the partition function ( $\approx$  the sequence space) induced by the relaxation. The preprocessing step would retain the same complexity, but based on a (reduced) treewidth  $tw' \leq tw$  for the relaxed graph  $G' = ([1, n], S' \cup \mathcal{D}')$ .

These complexities enable a tradeoff between the rejection (time), and the preprocessing (space), which may be critical to unlock future applications of RNA design. Indeed, the treewidth can be decreased by removing relatively few base pairs, as demonstrated below using our algorithm on pairs inferred for hard design instances.

We considered sets of DBPs inferred by RNAPond over two puzzles in the EteRNA benchmark. The EteRNA22 puzzle is an empty secondary structure spanning 400 nts, for which RNAPond obtains a valid design after inferring 465 DBPs. A tree decomposition of the graph formed by these 465 DPBs is then obtained with the standard min-fill-ordering heuristic [8], giving a width of 6. The EteRNA77 puzzle is 105 nts long, and consists in a collection of helices interspersed with destabilizing internal loops. RNAPond failed to produce a solution, and its final set of DBPs consists of 183 pairs, for which the same heuristic yields a tree decomposition of width 9. We further make both tree decompositions binary through bag duplications (see Appendix A), giving an FPT runtime to our algorithm, while potentially lowering the number of lost edges.

Executing the tree-diet algorithm (Theorem 10) on both graphs and their tree decompositions, we obtained simplified graphs, having lower treewidth while typically losing few edges, as illustrated and reported in Figure 7. Remarkably, the treewidth of the DBPs inferred for EteRNA22 can be decreased to  $tw' = 5$  by only removing 5 DBPs/edges (460/465 retained), and to  $tw' = 4$  by removing 4 further DBPs (456/465). For EteRNA77, our algorithm reduces the treewidth from 9 to 6 by only removing 7 DBPs.

Rough estimates can be provided for the tradeoff between the rejection and preprocessing complexities, by assuming that removing a DBP homogeneously increases the value of  $\mathcal{Z}$  by a factor  $\alpha := 16/10$  ( $\#pairs/\#incomp. pairs$ ). The relative increase in partition function is then  $\bar{q} \approx \alpha^b$ , when  $b$  base pairs are removed. For EteRNA22, reducing the treewidth by 2 units ( $6 \rightarrow 4$ ), *i.e.* a 16 fold reduction of the memory and preprocessing time, can be achieved by removing 9 DBPs, *i.e.* a 69 fold expected increase in the time of the generation phase. For EteRNA77, the same 16 fold ( $tw' = 9 \rightarrow 7$ ) reduction of the preprocessing time/space can be achieved through an estimated 4 fold increase of the generation time. A more aggressive 256 fold memory gain can be achieved at the expense of an estimated 1 152 fold increase in generation time. Given the large typical asymmetry in runtimes and implementation constants between the computation-heavy preprocessing and, relatively light, generation phases, the availability of an algorithm for the TREE-DIET problem provides new options, especially to circumvent memory limitations.

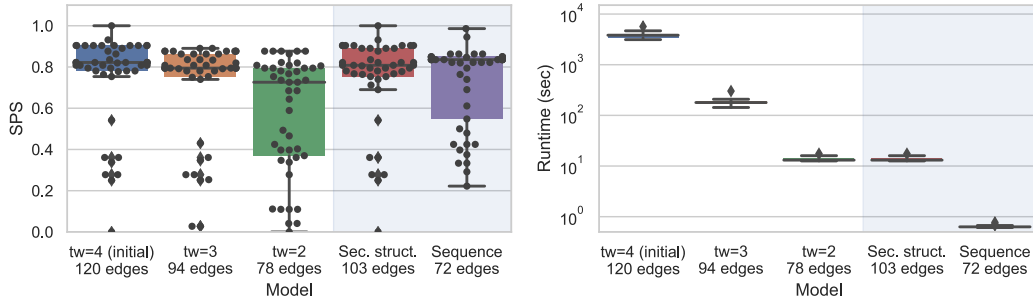
## 5.2 Structural alignment of complex RNAs

Structural homology is often posited within functional families of non-coding RNAs, and is foundational to algorithmic methods for multiple RNA alignments [23], considering RNA base pairs while aligning distant homologs. In the presence of complex structural features (pseudoknots, base triplets), the sequence-structure alignment problem becomes hard, yet admits XP solutions based on the treewidth of the base pair + backbone graph. In particular, Rinaudo *et al.* [32] describe a  $\Theta(n.m^{tw+1})$  algorithm for optimally aligning a structured RNA of length  $n$  onto a genomic region of length  $m$ . It optimizes an alignment score that includes: i) substitution costs for matches/mismatches of individual nucleotides and base pairs (including arc-breaking) based on the RIBOSUM matrices [24]; and ii) an affine gap cost model [33]. We used the implementation of the Rinaudo *et al.* algorithm, implemented in the LicoRNA software package [45, 46].

### 5.2.1 Impact of treewidth on the structural alignment of a riboswitch

In this case study, we used our tree-diet algorithm to modulate the treewidth of complex RNA structures, and investigate the effect of the simplification on the quality and runtimes of structure-sequence alignments. We considered the Cyclic di-GMP-II riboswitch, a regulatory motif found in bacteria that is involved in signal transduction, and undergoes conformational change upon binding the second messenger c-di-GMP-II [40, 41]. A 2.5Å resolution 3D model of the c-di-GMP-II riboswitch in *C. acetobutylicum*, proposed by Smith *et al.* [38] based on X-ray crystallography, was retrieved from the PDB [4] (PDBID: 3Q3Z). We annotated its base pairs geometrically using the DSSR method [28]. The canonical base pairs, supplemented with the backbone connections, were then accumulated in a graph, for which we heuristically computed an initial tree decomposition  $\mathcal{T}_4$ , having treewidth  $tw = 4$ .

We simplified our the initial tree decomposition  $\mathcal{T}_4$ , and obtained simplified models  $\mathcal{T}_3$ , and  $\mathcal{T}_2$ , having width  $tw' = 3$  and 2 respectively. As controls, we included tree decompositions based on the secondary structure (max. non-crossing set of BPs;  $\mathcal{T}_{2D}$ ) and sequence ( $\mathcal{T}_{1D}$ ).



■ **Figure 8** Impact on alignment quality (SPS; Left) and runtime (Right) of simplified instances for the RNA sequence-structure alignment of the pseudoknotted c-di-GMP-II riboswitch. The impact of simplifications on the quality of predicted alignments, using RFAM RF01786 as a reference, appears limited while the runtime improvement is substantial.

We used LicoRNA to predict an alignment  $a_{\mathcal{T},w}$  of each original/simplified tree decomposition  $\mathcal{T}$  onto each sequence  $w$  of the c-di-GMP-II riboswitch family in the RFAM database [23] (RF01786). Finally, we reported the LicoRNA runtime, and computed the Sum of Pairs Score (SPS) [43] as a measure of the accuracy of  $a_{\mathcal{T},w}$  against a reference alignment  $a_w^*$ :

$$\text{SPS}(a_{\mathcal{T},w}; a_w^*) = \frac{|\text{MatchedCols}(a_{\mathcal{T},w}) \cap \text{MatchedCols}(a_w^*)|}{|\text{MatchedCols}(a_w^*)|},$$

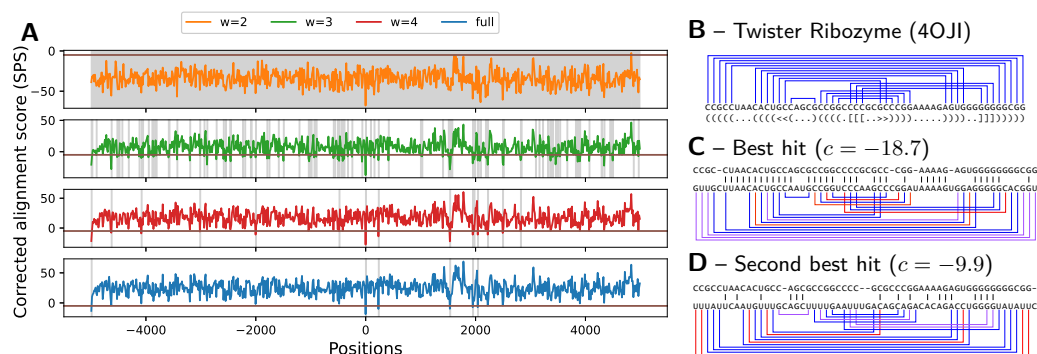
using as reference the alignment  $a_w^*$  between the 3Q3Z sequence and  $w$  induced by the manually-curated RFAM alignment of the RF01786 family.

The results, presented in Figure 8, show a limited impact of the simplification on the quality of the predicted alignment, as measured by the SPS in comparison with the RFAM alignment. The best average SPS (77.3%) is achieved by the initial model, having treewidth of 4, but the average difference with simplified models appears very limited (e.g. 76.5% for  $\mathcal{T}_3$ ), especially when considering the median. Meanwhile, the runtimes mainly depend on the treewidth, ranging from 1h for  $\mathcal{T}_4$  to 300ms for  $\mathcal{T}_{1D}$ . Overall,  $\mathcal{T}_{2D}$  seems to represent the best compromise between runtime and SPS, although its SPS may be artificially inflated by our election of RF01786 as our reference (built from a covariance model, *i.e.* essentially a 2D structure). Finally, the difference in number of edges (and induced SPS) between  $\mathcal{T}_{2D}$  and  $\mathcal{T}_2$ , both having  $tw = 2$ , exemplifies the difference between the TREE-DIET and GRAPH-DIET problems, and motivates further work on the latter.

## 5.2.2 Exact iterative strategy for the genomic search of ncRNAs

In this final case study, we consider an exact filtering strategy to search new occurrences of a structured RNA within a given genomic context. In this setting, one attempts to find all  $\varepsilon$ -admissible (cost  $\leq \varepsilon$ ) occurrences/hits of a structured RNA  $S$  of length  $n$  within a given genome of length  $g \gg n$ , broken down in windows of length  $\kappa \cdot n$ ,  $\kappa > 1$ . Classically, one would align  $S$  against individual windows, and report those associated with an admissible alignment cost. This strategy would have an overall  $\Theta(g \cdot n^{tw+2})$  time complexity, applying for instance the algorithm of [32].

Our instance simplification framework enables an alternative strategy, that incrementally filters out unsuitable windows based on models of increasing granularity. Indeed, for any given target sequence, the min alignment cost  $c_\delta$  obtained for a simplified instance of treewidth  $tw - \delta$  can be corrected (*cf* Appendix D) into a lower bound  $c_\delta^*$  for the min alignment



**Figure 9** Corrected costs associated with the search for structured homologs of the Twister ribozyme in chromosome 5 of *S. bicolor*, using simplified instances of various treewidth (A). Gray areas represent scores which, upon correction, remain below the cutoff, and have to be considered for further steps of the iterated filtering. Canonical base pairs of the ribozyme (PDBID 4OJI; B), mapped onto to the best hit (C) and second best hit (D) found along the search colored depending on their support in the target sequence (Red: incompatible; Purple: unstable G-U; Blue: stable).

cost  $c_0^*$  of the full-treewidth instance  $tw$ . Any window such that  $c_3^* > \varepsilon$  thus also obeys  $c_0^* > \varepsilon$ , and can be safely discarded from the list of putative  $\varepsilon$ -admissible windows, without having to perform a full-treewidth alignment. Given the exponential growth of the alignment runtime for increasing treewidth values (see Figure 8-right) this strategy is expected to yield substantial runtime savings.

We used this strategy to search occurrences of the Twister ribozyme (PDBID 4OJI), a highly-structured ( $tw = 5$ ) 54nts RNA initially found in *O. sativa* (asian rice) [26]. We targeted the *S. bicolor* genome (sorghum), focusing on a 10kb region centered on the 2,485,140 position of the 5th chromosome, where an instance of the ribozyme was suspected within an uncharacterized transcript (LOC110435504). The 4OJI sequence and structure were extracted from the 3D model as above, and included into a tree decomposition  $\mathcal{T}_5$  (73 edges), simplified into  $\mathcal{T}_4$  (71 edges),  $\mathcal{T}_3$  (68 edges) and  $\mathcal{T}_2$  (61 edges) using the tree-diet algorithm.

We aligned all tree decompositions against all windows of size 58nts using 13nts offset, and measured the score and runtime of the iterative filtering strategy using a cost cutoff  $\varepsilon = -5$ . The search recovers the suspected occurrence of twister as its best result (Figure 9.C), but produced hits (*cf* Figure 9.D) with comparable sequence conservation that could be the object of further studies. Regarding the filtering strategy, while  $\mathcal{T}_2$  only allows to rule out 3 windows out of 769,  $\mathcal{T}_3$  allows to eliminate an important proportion of putative targets, retaining only 109 windows, further reduced to 15 windows by  $\mathcal{T}_4$ , 6 of which end up as final hits for the full model  $\mathcal{T}_5$  (*cf* Figure 9.A). The search remains exact, but greatly reduces the overall runtime from 24 hours to 34 minutes (42 fold!).

## 6 Conclusion and discussion

We have established the parameterized complexity of three treewidth reduction problems, motivated by applications in Bioinformatics, as well as proposed practical algorithms for instances of reasonable treewidths. The reduced widths obtained by our proposed algorithm can be used to obtain: i) sensitive heuristics, owing to the consideration of a maximal amount of edges/information in the thinned graphs; ii) *a posteriori* approximation ratios, by comparing the potential contribution of removed edges to the optimal score obtained of the

thinned instance by a downstream FPT/XP algorithm; iii) substantial practical speedups without loss of correctness, e.g. when partial filtering can be safely achieved based on simplified input graphs

**Open questions.** Regarding the parameterized complexity of GRAPH-DIET and TREE-DIET, some questions remain open (see Table 1): an FPT algorithm for TREE-DIET (ideally, with  $2^{O(tw)} \cdot n$  running time), would be the most desirable, if possible satisfying the backbone constraints. We also aim at settling the parameterized complexity of the GRAPH-DIET problem, and try to give efficient exact algorithms for this problem (possibly using some tree-decomposition in input). Finally, we did not include the number of deleted edges in our multivariate analysis: even though in practice it is more difficult, a priori, to guarantee it has a small value, we expect it can be used to improve the running time in many cases.

**Backbone Preservation.** In two of our applications, the RNA secondary structure graph contains two types of edges: those representing the *backbone* of the sequence (i.e., between consecutive bases) and those representing base pair bounds. In practice, we want all backbone edges to be visible in the resulting tree-decomposition, and only base pairs may be lost. This can be integrated to the TREE-DIET model (and to our algorithms) using weighted edges, using the total weight rather than the count of deleted edges for the objective function. Note that some instances might be unrealizable (with no tree-diet preserving the backbone, especially for low  $tw'$ ). In most cases, ad-hoc bag duplications can help avoid this issue.

From a theoretical perspective, weighted edges may only increase the algorithmic complexity of the problems. However, a more precise model could consider graphs which already include a hamiltonian path (the backbone), and the remaining edges form a degree-one or two subgraph. Such extra properties may, in some cases, actually reduce the complexity of the problem. As an extreme case, we conjecture the PATH-DIET problem for  $tw' = 1$  becomes polynomial in this setting.

---

## References

- 1 Tatsuya Akutsu. Dynamic programming algorithms for RNA secondary structure prediction with pseudoknots. *Discrete Appl. Math.*, 104(1-3):45–62, 2000. doi:10.1016/S0166-218X(00)00186-4.
- 2 Julien Baste, Christophe Paul, Ignasi Sau, and Celine Scornavacca. Efficient FPT algorithms for (strict) compatibility of unrooted phylogenetic trees. *Bulletin of Mathematical Biology*, 79(4):920–938, February 2017. doi:10.1007/s11538-017-0260-y.
- 3 Julien Baste, Ignasi Sau, and Dimitrios M. Thilikos. Hitting minors on bounded treewidth graphs. i. general upper bounds. *SIAM J. Discret. Math.*, 34(3):1623–1648, 2020. doi:10.1137/19M1287146.
- 4 H. M. Berman, J. Westbrook, Z. Feng, G. Gilliland, T. N. Bhat, H. Weissig, I. N. Shindyalov, and P. E. Bourne. The protein data bank. *Nucleic acids research*, 28:235–242, 2000. doi:10.1093/nar/28.1.235.
- 5 Guillaume Blin, Alain Denise, Serge Dulucq, Claire Herrbach, and Helene Touzet. Alignments of RNA structures. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 7(2):309–322, April 2010. doi:10.1109/tcbb.2008.28.
- 6 Hans L Bodlaender. A linear-time algorithm for finding tree-decompositions of small treewidth. *SIAM Journal on computing*, 25(6):1305–1317, 1996.
- 7 Hans L Bodlaender and Arie MCA Koster. Combinatorial optimization on graphs of bounded treewidth. *The Computer Journal*, 51(3):255–269, 2008.



- 8 Hans L Bodlaender and Arie MCA Koster. Treewidth computations i. upper bounds. *Information and Computation*, 208(3):259–275, 2010.
- 9 Laurent Bulteau, Guillaume Fertin, Minghui Jiang, and Irena Rusu. Tractability and approximability of maximal strip recovery. *Theoretical Computer Science*, 440:14–28, 2012.
- 10 Laurent Bulteau and Mathias Weller. Parameterized algorithms in bioinformatics: An overview. *Algorithms*, 12(12):256, December 2019. doi:10.3390/a12120256.
- 11 Leizhen Cai. Fixed-parameter tractability of graph modification problems for hereditary properties. *Information Processing Letters*, 58(4):171–176, 1996.
- 12 Leizhen Cai. Parameterized complexity of vertex colouring. *Discrete Applied Mathematics*, 127(3):415–429, 2003.
- 13 Bruno Courcelle. The monadic second-order logic of graphs. i. recognizable sets of finite graphs. *Information and Computation*, 85(1):12–75, 1990. doi:10.1016/0890-5401(90)90043-H.
- 14 Christophe Crespelle, Pål Grønås Drange, Fedor V Fomin, and Petr A Golovach. A survey of parameterized algorithms and the complexity of edge modification. *arXiv preprint*, 2020. arXiv:2001.06867.
- 15 Marek Cygan, Fedor V Fomin, Łukasz Kowalik, Daniel Lokshtanov, Dániel Marx, Marcin Pilipczuk, Michał Pilipczuk, and Saket Saurabh. *Parameterized algorithms*, volume 5. Springer, 2015.
- 16 Marek Cygan, Daniel Lokshtanov, Marcin Pilipczuk, Michał Pilipczuk, and Saket Saurabh. On the hardness of losing width. In *International Symposium on Parameterized and Exact Computation*, pages 159–168. Springer, 2011.
- 17 Rodney G Downey and Michael Ralph Fellows. *Parameterized complexity*. Springer Science & Business Media, 2012.
- 18 Ehab S El-Mallah and Charles J Colbourn. The complexity of some edge deletion problems. *IEEE transactions on circuits and systems*, 35(3):354–362, 1988.
- 19 Vibhav Gogate and Rina Dechter. A complete anytime algorithm for treewidth. *arXiv preprint*, 2012. arXiv:1207.4109.
- 20 Stefan Hammer, Wei Wang, Sebastian Will, and Yann Ponty. Fixed-parameter tractable sampling for RNA design with multiple target structures. *BMC Bioinformatics*, 20(1), April 2019. doi:10.1186/s12859-019-2784-7.
- 21 Buhm Han, Banu Dost, Vineet Bafna, and Shaojie Zhang. Structural alignment of pseudoknotted RNA. *Journal of Computational Biology*, 15(5):489–504, 2008. doi:10.1089/cmb.2007.0214.
- 22 Wenzel Jakob, Jason Rhinelander, and Dean Moldovan. pybind11 – seamless operability between c++11 and python, 2017. URL: <https://github.com/pybind/pybind11>.
- 23 Ioanna Kalvari, Eric P Nawrocki, Nancy Ontiveros-Palacios, Joanna Argasinska, Kevin Lamkiewicz, Manja Marz, Sam Griffiths-Jones, Claire Toffano-Nioche, Daniel Gautheret, Zasha Weinberg, Elena Rivas, Sean R Eddy, Robert D Finn, Alex Bateman, and Anton I Petrov. Rfam 14: expanded coverage of metagenomic, viral and microRNA families. *Nucleic Acids Research*, 49(D1):D192–D200, November 2020. doi:10.1093/nar/gkaa1047.
- 24 Robert J Klein and Sean R Eddy. Rsearch: finding homologs of single structured RNA sequences. *BMC bioinformatics*, 4(1):44, 2003.
- 25 Neocles B Leontis and Eric Westhof. Geometric nomenclature and classification of RNA base pairs. *RNA*, 7(4):499–512, 2001.
- 26 Yijin Liu, Timothy J Wilson, Scott A McPhee, and David MJ Lilley. Crystal structure and mechanistic investigation of the twister ribozyme. *Nature chemical biology*, 10(9):739–744, 2014.
- 27 László Lovász. Graph minor theory. *Bulletin of the American Mathematical Society*, 43(1):75–86, 2006.
- 28 Xiang-Jun Lu, Harmen J. Bussemaker, and Wilma K. Olson. DSSR: an integrated software tool for dissecting the spatial structure of RNA. *Nucleic Acids Research*, 43(21):e142–e142, July 2015. doi:10.1093/nar/gkv716.

- 29 R. B. Lyngsø and C. N. S. Pedersen. RNA pseudoknot prediction in energy-based models. *Journal of Computational Biology*, 7(3-4):409–427, 2000.
- 30 Andrzej Proskurowski and Jan Arne Telle. Classes of graphs with restricted interval models. *Discrete Mathematics & Theoretical Computer Science*, 3(4), 2006.
- 31 Vladimir Reinharz, Antoine Soulé, Eric Westhof, Jérôme Waldispühl, and Alain Denise. Mining for recurrent long-range interactions in RNA structures reveals embedded hierarchies in network families. *Nucleic Acids Research*, 46(8):3841–3851, March 2018. doi:10.1093/nar/gky197.
- 32 Philippe Rinaudo, Yann Ponty, Dominique Barth, and Alain Denise. Tree decomposition and parameterized algorithms for RNA structure-sequence alignment including tertiary interactions and pseudoknots. In *Lecture Notes in Computer Science*, pages 149–164. Springer Berlin Heidelberg, 2012. doi:10.1007/978-3-642-33122-0\_12.
- 33 Elena Rivas and Sean R Eddy. Parameterizing sequence alignment with an explicit evolutionary model. *BMC bioinformatics*, 16(1):406, 2015.
- 34 Neil Robertson and Paul D Seymour. Graph minors. xiii. the disjoint paths problem. *Journal of combinatorial theory, Series B*, 63(1):65–110, 1995.
- 35 Toshiki Saitoh, Ryo Yoshinaka, and Hans L. Bodlaender. Fixed-treewidth-efficient algorithms for edge-deletion to interval graph classes. In *WALCOM: Algorithms and Computation - 15th International Conference and Workshops, 2021, Proceedings*, volume 12635 of *Lecture Notes in Computer Science*, pages 142–153. Springer, 2021. doi:10.1007/978-3-030-68211-8\_12.
- 36 Roman Sarrazin-Gendron, Hua-Ting Yao, Vladimir Reinharz, Carlos G. Oliver, Yann Ponty, and Jérôme Waldispühl. Stochastic sampling of structural contexts improves the scalability and accuracy of RNA 3d module identification. In *Lecture Notes in Computer Science*, pages 186–201. Springer International Publishing, 2020. doi:10.1007/978-3-030-45257-5\_12.
- 37 Saad Sheikh, Rolf Backofen, and Yann Ponty. Impact Of The Energy Model On The Complexity Of RNA Folding With Pseudoknots. In Juha Kärkkäinen and Jens Stoye, editors, *CPM - 23rd Annual Symposium on Combinatorial Pattern Matching - 2012*, volume 7354 of *Combinatorial Pattern Matching*, pages 321–333, Helsinki, Finland, 2012. Juha Kärkkäinen, Springer. doi:10.1007/978-3-642-31265-6\_26.
- 38 Kathryn D. Smith, Carly A. Shanahan, Emily L. Moore, Aline C. Simon, and Scott A. Strobel. Structural basis of differential ligand recognition by two classes of bis-(3′-5′)-cyclic dimeric guanosine monophosphate-binding riboswitches. *Proceedings of the National Academy of Sciences*, 108(19):7757–7762, 2011. doi:10.1073/pnas.1018857108.
- 39 Yinglei Song, Chunmei Liu, Russell Malmberg, Fangfang Pan, and Liming Cai. Tree decomposition based fast search of RNA structures including pseudoknots in genomes. In *Computational Systems Bioinformatics Conference, 2005. Proceedings. 2005 IEEE*, pages 223–234. IEEE, 2005.
- 40 N. Sudarsan, E. R. Lee, Z. Weinberg, R. H. Moy, J. N. Kim, K. H. Link, and R. R. Breaker. Riboswitches in eubacteria sense the second messenger cyclic di-gmp. *Science*, 321(5887):411–413, 2008. doi:10.1126/science.1159519.
- 41 Rita Tamayo. Cyclic diguanylate riboswitches control bacterial pathogenesis mechanisms. *PLOS Pathogens*, 15(2):1–7, February 2019. doi:10.1371/journal.ppat.1007529.
- 42 Jinsong Tan and Louxin Zhang. The consecutive ones submatrix problem for sparse matrices. *Algorithmica*, 48(3):287–299, 2007.
- 43 J D Thompson, F Plewniak, and O Poch. BALiBASE: a benchmark alignment database for the evaluation of multiple alignment programs. *Bioinformatics*, 15(1):87–88, January 1999. doi:10.1093/bioinformatics/15.1.87.
- 44 Jelena Vucinic, David Simoncini, Manon Ruffini, Sophie Barbe, and Thomas Schiex. Positive multistate protein design. *Bioinformatics*, 36(1):122–130, June 2019. doi:10.1093/bioinformatics/btz497.
- 45 Wei Wang. *Practical sequence-structure alignment of RNAs with pseudoknots*. PhD thesis, Université Paris-Saclay, School of Computer Science, 2017.

- 46 Wei Wang, Alain Denise, and Yann Ponty. Licorna: alignment of complex rnas v1.0, 2017. URL: <https://licorna.lri.fr>.
- 47 M. S. Waterman. Secondary structure of single stranded nucleic acids. *Advances in Mathematics Supplementary Studies*, 1(1):167–212, 1978.
- 48 Mathias Weller, Annie Chateau, and Rodolphe Giroudeau. Exact approaches for scaffolding. *BMC Bioinformatics*, 16(S14), October 2015. doi:10.1186/1471-2105-16-s14-s2.
- 49 A. Xayaphoummine, T. Bucher, F. Thalmann, and H. Isambert. Prediction and statistics of pseudoknots in RNA structures using exactly clustered stochastic simulations. *Proc. Natl. Acad. Sci. U. S. A.*, 100(26):15310–15315, 2003.
- 50 Jinbo Xu. Rapid protein side-chain packing via tree decomposition. In *Lecture Notes in Computer Science*, pages 423–439. Springer Berlin Heidelberg, 2005. doi:10.1007/11415770\_32.
- 51 Hua-Ting Yao, Jérôme Waldspühl, Yann Ponty, and Sebastian Will. Taming Disruptive Base Pairs to Reconcile Positive and Negative Structural Design of RNA. In *RECOMB 2021 - 25th international conference on research in computational molecular biology*, Padova, France, 2021.

## A Editing Trees before the Diet

Any tree decomposition can be transformed into a binary one through the duplications of bags having more than 2 children. To do so in practice, one will, as long as the tree decomposition is not binary, apply the following transformation:

1. Find a bag  $X$  with children  $Y_1, \dots, Y_\Delta$  and  $\Delta > 2$ .
2. Introduce a new bag  $X'$  with the same content as  $X$  and locally modify the tree decomposition in the following way:  $X$  will now have  $Y_1$  and  $X'$  as children, while  $X'$  will have  $Y_2 \cdots Y_\Delta$ .

When it is no longer possible to apply this transformation, the tree decomposition is binary. For each bag having originally  $\Delta > 2$  children in the decomposition,  $\Delta - 1$  new bags have been introduced. In total, with  $N_{bags}$  the original number of bags in the decomposition, strictly less than  $N_{bags}$  new bags have been introduced (each new bag is associated to an edge of the original tree decomposition).

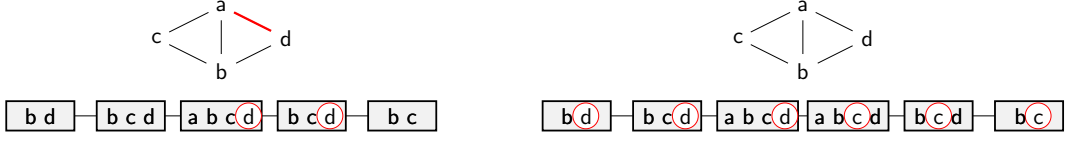
This transformation is in fact the first step towards obtaining a *nice tree decomposition* [7, 15].

A question that arises then is what impact these modifications may have on the output of TREE-DIET, when applied to the tree decomposition given as input. We argue that duplication operations (as used above to get a binary tree decomposition) can only improve the solution, i.e decrease the number of *lost* edges. Indeed, within the coloring formulation of the problem, new bags yield new opportunities for an edge to be *represented*, with both its end-points green in some bag. See Figure 10 for an illustration.

More generally, any operation on the input tree decomposition that does not suppress any of the original bags can only improve the solution to the TREE DIET problem. We do not tackle here the problem of finding the best edition operations to apply onto a tree decomposition given as input to TREE DIET, which is an a priori difficult task.

## B Pseudo-code

Algorithm 1 present a pseudo-code of our dynamic programming algorithm for TREE DIET, with a memoization approach. The C++/pybind11 [22] implementation is available at <https://gitlab.inria.fr/amibio/tree-diet>.



■ **Figure 10** Left: A graph and a path-decomposition whose optimal 1-tree-diet loses an edge ( $ad$ ). However, duplicating the bag  $abcd$  (right) yields a tree-decomposition with a lossless 1-tree-diet.

Note that the implementation allows to solve a more general *weighted* version of TREE DIET, where each edge is given a weight, and the objective is to find a  $(tw - tw')$ -diet of the input tree decomposition preserving a set of edges of maximum total weight.

In the context of RNA applications, this feature allows to favour as much as possible preservation of the backbone of RNA molecules, i.e. edges between consecutive nucleotides along the string, by assigning them a weight greater than the number of non-backbone edges.

Edge weights are passed to the function in the form of a dictionary/map  $W$  associating a real weight to each edge. Within Algorithms 1, the only place where it is taken into account is the the *count* function, which computes the weight of edges accounted for by the bag that is currently visited.

## C Correctness of the rejection-based sampling of RNA designs

A recent method for RNA design, called RNAPond [51], implements a sampling approach to tackle the inverse folding of RNA. Targeting a secondary structure  $S$  of length  $n$ , it performs a Boltzmann-weighted sampling of sequences and, at each iteration, identifies Disruptive Base Pairs (DBPs) that are not in  $S$ , yet are recurrent in the Boltzmann ensemble of generated sequences. Those base pairs are then added to a set  $\mathcal{D}$  of DBPs, and excluded in subsequent generations through an assignment of non-binding pairs of nucleotides, outside of  $\mathcal{B} := \{(G, C), (C, G), (A, U), (U, A), (G, U), (U, G)\}$ .

At the core of the method, one finds a random generation algorithm which takes as input a secondary structure  $S$  and a set  $\mathcal{D}$  of DBPs. The algorithm generates from the set  $\mathcal{W}_{S,\mathcal{D}}$  of sequences  $w \in \{A, C, G, U\}^n$  which are: i) compatible with all  $(i, j) \in S$ , i.e.  $(w_i, w_j) \in \mathcal{B}$ ; and ii) incompatible with all  $(k, l) \in \mathcal{D}$ , i.e.  $(w_k, w_l) \notin \mathcal{B}$ . The algorithm then enforces a (dual) Boltzmann distribution over the sequences in  $\mathcal{W}_{S,\mathcal{D}}$ :

$$\forall w \in \mathcal{W}_{S,\mathcal{D}} : \mathbb{P}(w \mid \mathcal{D}, S) = \frac{e^{-\beta \cdot E_{w,S}}}{\mathcal{Z}_{S,\mathcal{D}}} \quad \text{with} \quad \mathcal{Z}_{S,\mathcal{D}} := \sum_{w' \in \mathcal{W}_{S,\mathcal{D}}} e^{-\beta \cdot E_{w',S}} \quad (1)$$

where  $\beta > 0$  is an arbitrary constant akin to a temperature. Yao *et al.* describe an algorithm which generates  $k$  sequences in  $\Theta(k(n + |\mathcal{D}|))$  time, after a preprocessing in  $\Theta(n \cdot |\mathcal{D}| \cdot 4^{tw})$  time and  $\Theta(n \cdot 4^{tw})$  space, where  $tw$  is the treewidth of the graph having edges in  $S \cup \mathcal{D}$ .

The discrepancy in the preprocessing and sampling complexities suggests an alternative strategy, utilizing rejection on top of a relaxed sampling. Namely, we consider a rejection algorithm, which starts from a relaxation  $(S', \mathcal{D}')$  of the initial constraints  $(S' \cup \mathcal{D}' \subset S \cup \mathcal{D})$ , and iterates Yao *et al.*'s algorithm to generate sequences in  $\mathcal{W}_{S',\mathcal{D}'} \supset \mathcal{W}_{S,\mathcal{D}}$ , rejecting those outside of  $\mathcal{W}_{S,\mathcal{D}}$ , until  $k$  suitable ones are obtained. The rejection algorithm generates a given sequence  $w \in \mathcal{W}_{S,\mathcal{D}}$  on its first attempt with probability  $p := e^{-\beta \cdot E_{w,S}} / \mathcal{Z}_{S',\mathcal{D}'}$  and, more generally, after  $r$  rejections with probability  $(1 - q)^r p$  with  $q := \mathcal{Z}_{S,\mathcal{D}} / \mathcal{Z}_{S',\mathcal{D}'}$ . The overall probability of emitting  $w$  is thus

$$p \cdot \sum_{r \geq 0} (1 - q)^r = \frac{p}{q} = \frac{e^{-\beta \cdot E_{w,S}}}{\mathcal{Z}_{S,\mathcal{D}}} = \mathbb{P}(w \mid \mathcal{D}, S).$$

■ **Algorithm 1** Dynamic programming algorithm for TREE-DIET.

---

**Input** : Tree-decomposition  $\mathcal{T}$ , graph  $G$ , target width  $tw'$ , edge weights  $W$   
**Output** : Maximum total weight of a set of realizable/non-lost edges in a  $(tw - tw')$ -diet of  $\mathcal{T}$   
**Side-Product** : A filled table  $c[X_i, f]$ ,  $\forall X_i$  bag and  $f$  coloring of  $X_i$

```

1 Function optim_num_real_edges( $X_i, f, G, tw', W$ ):
2   if  $c[X_i, f]$  already computed then return  $c[X_i, f]$ ; ;
3   if  $|f^{-1}(o) \cup f^{-1}(r)| \leq (|X_i| - tw' - 1)$  then
4     //not enough removals.;
5      $c[X_i, f] = -\infty$ ;
6     return  $c[X_i, f]$ ;
7   end
8   if  $X_i == \text{leaf}$  then
9      $c[X_i, f] = 0$ ;
10    return  $c[X_i, f]$ ;
11  end
12  int ans =  $-\infty$ ;
13  for  $m \in \text{orange\_maps}(X_i, f)$  do
14    int ans_m = 0;
15    for  $Y_j \in X_i.\text{children}$  do
16      int ans_j =  $-\infty$ ;
17      for  $f'_j \in \text{compatible}(f, m, X_i, Y_j)$  do
18        int val = 0;
19        val +=  $\text{count}(f, f'_j, W)$ ;
20        val += optim_num_real_edges( $Y_j, f'_j, G, tw'$ );
21        if val ≥ ans_j then ans_j = val;
22      ;
23    end
24    ans_m += ans_j;
25  end
26  if ans_m ≥ ans then ans = ans_m; ;
27 end
28  $c[X_i, f] = \text{ans}$ 
29 return  $c[X_i, f]$ ;
30 end

```

---

In other words, our relaxed generator coupled with the rejection step, represents an unbiased algorithm for the Boltzmann distribution of Eq. (1) over  $\mathcal{W}_{S, \mathcal{D}}$ .

Meanwhile, the average-case complexity can be impacted by the strategy. Indeed, the relaxed instance  $(S', \mathcal{D}')$  can accelerate the preprocessing due to a reduced treewidth  $tw' \leq tw$ . The rejection step only increases the expected number of generations by a factor  $\bar{q} := \mathcal{Z}_{S', \mathcal{D}'} / \mathcal{Z}_{S, \mathcal{D}}$ , representing the inflation of the sequence space, induced by the relaxation of the constraints. Overall, the average-case time complexity of the rejection algorithm is in  $\Theta(n \cdot |\mathcal{D}'| \cdot 4^{tw'} + k \cdot \bar{q} \cdot (n + |\mathcal{D}'|))$  time and  $\Theta(n \cdot 4^{tw'})$  space. This space improvement is notable when  $tw' < tw$ , and could be key for the practical applicability of the method, especially given that memory represents the bottleneck of most treewidth-based DP algorithms.

## D Lower bound for the min. alignment cost from simplified models

Here, we justify the filtering strategy described in Section 5.2.2. Namely, we formally prove that, given a structured RNA  $S$  and a targeted genomic region  $w$ , a lower bound for the minimal alignment cost of  $S$  and  $w$  can be obtained from the minimal alignment cost of some  $S' \subseteq S$  and  $w$ . If this lower bound for  $S' \subseteq S$  is higher than the specified cutoff  $\varepsilon$ , then there is no need to align  $w$  to  $S$  the full model, as the resulting cost is guaranteed to stay above the selection cutoff  $\varepsilon$ .

Let  $S$  be an arc-annotated sequence of length  $m$  ( $S_i$  denotes the  $i$ th character of  $S$ ),  $w$  be a target (flat) sequence of length  $m$ , and  $\mu : [1, n] \rightarrow [1, m] \cup \{\perp\}$  represents an alignment<sup>2</sup>. We consider the following cost function, adapted from [32], which quantifies the quality of an alignment  $\mu$  for  $S$  and  $w$ :

$$C(S, w, \mu) = \sum_{\substack{i \text{ unpaired in } S, \\ k := \mu_i}} \gamma(S_i, w_k) + \sum_{\substack{(i,j) \in S, \\ (k,l) := (\mu_i, \mu_j)}} \phi(S_i, S_j, w_k, w_l) \\ + \sum_{g \in \text{gaps}(S)} \lambda_g(g) + \sum_{g \in \text{gaps}(w)} \lambda_T(g)$$

where

- $\gamma(a, b)$  returns the *substitution cost* which penalizes (mismatches) or rewards (matches) the substitution of  $a$  into  $b$  (set to 0 and handled in gaps if  $b = \perp$ );
- $\phi(a, b, c, d)$  return a *base pair substitution cost*, penalizing (arc breaking) or rewarding (conservation or compensatory mutations) the transformation of nucleotides  $(a, b)$  into nucleotides/gaps  $(c, d)$  (set to 0 and handled in gaps if  $(c, d) = (\perp, \perp)$ );
- $\lambda_S$  and  $\lambda_T$  penalize gaps introduced by  $\mu$  respectively in  $S$  and  $w$  (affine cost model).

Given this definition, consider a simplified model  $S' \subset S$ , associated with a minimal cost

$$c' := \min_{\mu} C(S, w, \mu)$$

and denote by  $c^*$  the minimal cost of the full model  $S$ , we have the following inequality.

► **Proposition 14.**

$$c' - \sum_{\substack{i \text{ unpaired in } S', \\ \text{paired in } S}} \max_b \gamma(S_i, b) + \sum_{(i,j) \in S \setminus S'} \min_{a,b} \phi(S_i, S_j, a, b) \leq c^* \quad (2)$$

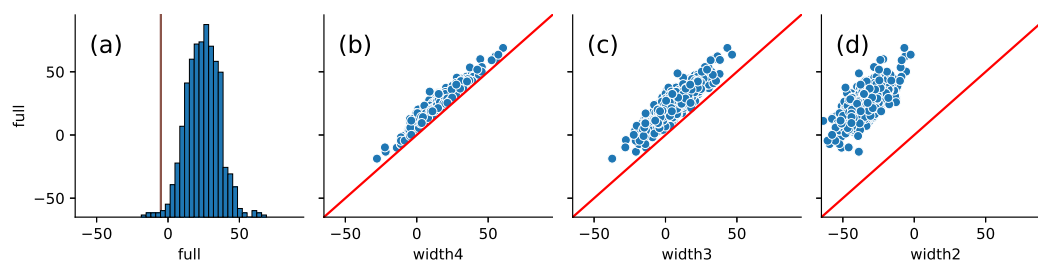
**Proof.** For any alignment, we have, per the definition of  $C(S, w, \mu)$ :

$$C(S, w, \mu) = C(S', w, \mu) - \sum_{\substack{i \text{ unpaired in } S', \\ \text{paired in } S, \\ \text{and } k := \mu_i}} \gamma(S_i, w_k) + \sum_{\substack{(i,j) \in S \setminus S' \\ \text{s.t. } (k,l) := (\mu_i, \mu_j)}} \phi(S_i, S_j, w_k, w_l).$$

Minimizing over all alignment  $\mu$ , one obtains

$$\min_{\mu} C(S, w, \mu) = \min_{\mu} C(S', w, \mu) - \sum_{\substack{i \text{ unpaired in } S', \\ \text{paired in } S, \\ \text{and } k := \mu_i}} \gamma(S_i, w_k) + \sum_{\substack{(i,j) \in S \setminus S' \\ \text{s.t. } (k,l) := (\mu_i, \mu_j)}} \phi(S_i, S_j, w_k, w_l).$$

<sup>2</sup> An alignment  $\mu$  is subject to further constraints, notably including some restricted form of monotonicity, when represented as a function. However, those constraints are reasonably intuitive and we omit them in this discussion for the sake of simplicity.



■ **Figure 11** (a) Histogram of alignment scores obtained by aligning the full structure ( $tw = 5$ ) model of the Twister ribozyme (pdb-id: 4OJI) with  $\kappa \cdot n$ -sized windows in a 10kb region of the 5<sup>th</sup> chromosome of *S. bicolor*. A vertical line is positioned at the  $\epsilon$  threshold. (b;c;d) Corrected alignment scores obtained for reduced-treewidth models for each window, plotted against the corresponding score of the full model. The corrected alignment score indeed acts as a lower bound to the full-model score (points above the  $y = x$  red line), allowing a iterative filtering strategy.

Independently minimizing each term of the right-hand-side, we obtain a first lower bound

$$c^* \geq c' - \max_{\mu} \sum_{\substack{i \text{ unpaired in } S', \\ \text{paired in } S, \\ \text{and } k:=\mu_i}} \gamma(S_i, w_k) + \min_{\mu} \sum_{\substack{(i,j) \in S \setminus S' \\ \text{s.t. } (k,l):=(\mu_i, \mu_j)}} \phi(S_i, S_j, w_k, w_l).$$

further coarsened by an independent optimization of the elements in the sums

$$\begin{aligned} c^* &\geq c' - \sum_{\substack{i \text{ unpaired in } S', \\ \text{paired in } S}} \max_{\mu} \gamma(S_i, w_k) + \sum_{(i,j) \in S \setminus S'} \min_{\mu} \phi(S_i, S_j, w_k, w_l) \\ &= c' - \sum_{\substack{i \text{ unpaired in } S', \\ \text{paired in } S}} \max_a \gamma(S_i, a) + \sum_{(i,j) \in S \setminus S'} \min_{a,b} \phi(S_i, S_j, a, b). \end{aligned}$$

where the last line is obtained by considering the worst-case contributors to nucleotides and base pairs substitutions. Importantly, the right-hand side no longer depends on  $\mu$  any more, and can be used to easily computed a corrected score/lower bound. ◀

The corrected expression, shown in the left hand side of Equation (2) allows, when lower than a cutoff  $\epsilon$ , to safely discard  $w$  as a potential hit for the full model  $S$ . This corrected score score is plotted in Figure 9A, allowing for a gradual reduction of the search space for  $\epsilon$ -admissible hits. We show in Figure 11 the corrected scores obtained for simplified structures  $S'$  of various treewidths, plotted against the scores of the full target structure.