



Made available through Montana State University's [ScholarWorks](#)

Computing the Tandem Duplication Distance is NP-Hard

Manuel Lafond, Binhai Zhu, Peng Zou

© SIAM

COMPUTING THE TANDEM DUPLICATION DISTANCE IS NP-HARD*

MANUEL LAFOND[†], BINHAI ZHU[‡], AND PENG ZOU[§]

Abstract. In computational biology, tandem duplication is an important biological phenomenon which can occur either at the genome or at the DNA level. A tandem duplication takes a copy of a genome segment and inserts it right after the segment — this can be represented as the string operation $AXB \Rightarrow AXXB$. Tandem exon duplications have been found in many species such as human, fly or worm, and have been largely studied in computational biology.

The *Tandem Duplication* (TD) distance problem we investigate in this paper is defined as follows: given two strings S and T over the same alphabet Σ , compute the smallest sequence of tandem duplications required to convert S to T . The natural question of whether the TD distance can be computed in polynomial time was posed in 2004 by Leupold et al. and had remained open, despite the fact that tandem duplications have received much attention ever since. In this paper, we focus on the special case when all characters of S are distinct. This is known as the *exemplar* TD distance, which is of special relevance in bioinformatics. We first prove that this problem is NP-hard when the alphabet size is unbounded, settling the 16-year old open problem. We then show how to adapt the proof to $|\Sigma| = 4$, hence proving the NP-hardness of the TD problem for any $|\Sigma| \geq 4$. One of the tools we develop for the reduction is a new problem called *Cost-Effective Subgraph*, for which we obtain W[1]-hardness results that might be of independent interest. We finally show that computing the exemplar TD distance between S and T is fixed-parameter tractable. Our results open the door to many other questions, and we conclude with several open problems.

Key words. Tandem duplication, Text processing, Formal languages, Computational genomics, FPT algorithm

AMS subject classifications. 68Q25, 68R10, 68W32

1. Introduction. Tandem duplication is a biological process that creates consecutive copies of a segment of a genome during DNA replication. Representing genomes as strings, this event transforms a string AXB into another string $AXXB$. This process is known to occur either at small scale at the nucleotide level, or at large scale at the genome level [5, 6, 7, 23, 30]. For instance, it is known that the Huntington disease is associated with the duplication of 3 nucleotides CAG [26], whereas at genome level, tandem duplications are known to involve multiple genes during cancer progression [27]. Furthermore, gene duplication is believed to be the main driving force behind evolution, and the majority of duplications affecting organisms are believed to be of the tandem type (see e.g. [31]). As a result, around 3% of the human genome are formed of tandem repeats.

For these reasons, tandem duplications have received significant attention in the last decades, both in practice and theory. The combinatorial aspects of tandem duplications have been studied extensively by computational biologists [2, 12, 14, 21, 32], one question of interest being to reconstruct the evolution of a cluster of tandem repeats by duplications that could have given rise to the observed sequences. In parallel, various formal language communities [9, 33, 25] have investigated the expressive

*Submitted to the editors DATE. An abstract of this paper was presented at STACS'20 [20], and was also presented as part of an invited talk at CSR'20 [34].

Funding: This work was funded by NSERC of Canada.

[†]Department of Computer Science, Universite de Sherbrooke, Sherbrooke, Quebec J1K 2R1, Canada (manuel.lafond@usherbrooke.ca).

[‡]Gianforte School of Computing, Montana State University, Bozeman, MT 59717-3880, USA (bhz@montana.edu).

[§]Gianforte School of Computing, Montana State University, Bozeman, MT 59717-3880, USA (peng.zou@student.montana.edu).

42 power of tandem duplications on strings.

43 From the latter perspective, a natural question arises: given a string S , what
 44 is the language that can be obtained starting from S and applying (any number
 45 of) tandem duplications, i.e. rules of the form $AXB \rightarrow AXXB$, where X can be any
 46 substring of S ? This question was first asked in 1984 in the context of so-called *copying*
 47 *systems* [11]. Combined with results from [3], it was shown that this language is
 48 regular if S is on a binary alphabet, but not regular for larger alphabets. These results
 49 were rediscovered 15 years later in [9, 33]. In [25], it was shown that the membership,
 50 inclusion and regularity testing problems on the language defined by S can all be
 51 decided in linear time (still on binary alphabets). In [25, 24, 17], similar problems are
 52 also considered on non-binary alphabets, when the length $|X|$ of duplicated strings is
 53 bounded by a constant. More recently, Cho et al. [8] introduced a tandem duplication
 54 system where the *depth* of a character, i.e. the number of “generations” it took to
 55 generate it, is considered. In [16, 18], the authors study the *expressive power* of
 56 tandem duplications, a notion based on the subsequences that can be obtained from
 57 various types of copying mechanisms.

58 More directly related to our work, Alon et al. [1] recently investigated the mini-
 59 mum number of duplications required to transform a string S into another string T .
 60 We call this the *Tandem Duplication (TD) distance*. More specifically, the authors
 61 show that on binary strings, the maximum TD distance between a square-free string
 62 S and a string T of length n is $\Theta(n)$. They also mention the unsolved algorithmic
 63 problem of computing the TD distance between S and T . In fact, this question was
 64 posed earlier in [25] (pp. 306, Open Problem 3) by Leupold et al. and has remained
 65 open ever since. We settle this open problem in this paper for an unbounded alpha-
 66 bet. (We then extend the result when the alphabet size is 4.) As will be seen, our
 67 technique is different from that used in [1], which only works for binary strings.

68 On the other hand, the TD distance is one of the many ways of comparing two
 69 genomes represented as strings in computational biology — other notable examples
 70 include breakpoint [15] and transpositions distances, the latter having recently been
 71 shown NP-hard in a celebrated paper of Bulteau et al. [4]. The TD distance has itself
 72 received special attention recently, owing to its role in cancer evolution [28].

73 **Our results.** In this paper, we solve the problem posed by Leupold et al. in 2004
 74 and show that computing the TD distance from a string S to a string T is NP-hard.
 75 We show that this result holds even if S is *exemplar*, i.e. if each character of S
 76 is distinct. Exemplar strings are commonly studied in computational biology [29], since
 77 they represent genomes that existed prior to duplication events. We note that simply
 78 deciding if S can be transformed into T by a sequence of TDs still has unknown
 79 complexity. In our case, we show that the hardness of minimizing TDs holds on
 80 instances in which such a sequence is guaranteed to exist.

81 As demonstrated by the transpositions distance in [4], obtaining NP-hardness re-
 82 sults for string distances can sometimes be an involving task. Our hardness reduction
 83 is also quite technical, and one of the tools we develop for it is a new problem we
 84 call the *Cost-Effective Subgraph*. In this problem, we are given a graph G with a
 85 cost c , and we must choose a subset X of $V(G)$. Each edge with both endpoints in
 86 X has a cost of $|X|$, every other edge costs c , and the goal is to find a subset X of
 87 minimum cost. We show that this problem is W[1]-hard for parameter $p + c$, where p
 88 is a parameter that asks if one can achieve a cost of at most $c|E(G)| - p$ (here $c|E(G)|$

89 is an upper bound on the cost¹). The problem enforces optimizing the tradeoff be-
 90 tween covering many edges versus having a large subset of high cost, which might be
 91 applicable to other problems. In our case it captures the main difficulty in computing
 92 TD distances. We then obtain some positive results by showing that if S is exemplar,
 93 then one can decide if S can be transformed into T using at most k duplications in
 94 time $2^{O(k^2)} + \text{poly}(n)$, where n is the length of T . The result is obtained through an
 95 exponential size kernel. All of our results concern strings with unbounded alphabet
 96 sizes. Finally, we conclude with several open problems that might be of interest to
 97 the theoretical computer science community.

98 This paper is organized as follows. In Section 2, we give basic definitions. In
 99 Section 3, we show that computing the TD distance is NP-hard through the Cost-
 100 Effective Subgraph problem when $|\Sigma|$ is infinite. In Section 4, we show that the same
 101 problem is NP-hard for any constant $|\Sigma| \geq 4$. In Section 5, we show that computing
 102 the exemplar TD distance is FPT. In Section 6, we conclude the paper with several
 103 open problems.

104 **2. Preliminary notions.** We borrow the string terminology and notation from
 105 [13]. In particular, $[n]$ denotes the set of integers $\{1, 2, \dots, n\}$. Unless stated otherwise,
 106 all the strings in the paper are on an alphabet denoted Σ . If S_1 and S_2 are two strings,
 107 we usually denote their concatenation by S_1S_2 . For a string S , we write $\Sigma(S)$ for the
 108 subset of characters of Σ that have at least one occurrence in S . A string S is called
 109 *exemplar* if $|S| = |\Sigma(S)|$, i.e. each character present in S occurs only once. A *substring*
 110 of S is a contiguous sequence of characters within S . A *prefix* (resp. *suffix*) of S is
 111 a substring that occurs at the beginning (resp. end) of S , i.e. if $S = S_1S_2$ for some
 112 strings S_1 and S_2 , then S_1 is a prefix of S and S_2 a suffix of S . A *subsequence* of S
 113 is a string that can be obtained by successively deleting characters from S .

114 A *tandem duplication* (TD) is an operation on a string S that copies a substring
 115 X of S and inserts the copy after the occurrence of X in S . In other words, a TD
 116 transforms $S = AXB$ into $AXXB$. Given another string T , we write $S \Rightarrow T$ if there
 117 exist strings A, B, X such that $S = AXB$ and $T = AXXB$. More generally, we write
 118 $S \Rightarrow_k T$ if there exist S_1, \dots, S_{k-1} such that $S \Rightarrow S_1 \Rightarrow \dots \Rightarrow S_{k-1} \Rightarrow T$. We also
 119 write $S \Rightarrow_* T$ if there exists some k such that $S \Rightarrow_k T$.

120 **DEFINITION 2.1.** *The TD distance $\text{dist}_{TD}(S, T)$ between two strings S and T*
 121 *is the minimum value of k satisfying $S \Rightarrow_k T$. If $S \Rightarrow_* T$ does not hold, then*
 122 *$\text{dist}_{TD}(S, T) = \infty$.*

123 We use the term *distance* here to refer to the number of TD operations from a
 124 string S to another string T , but one may note that *TD* is not a metric in the formal
 125 sense. In particular, dist_{TD} is not symmetric since duplications can only increase the
 126 length of a string.

127 A *square string* is a string of the form XX , i.e. a concatenation of two identical
 128 substrings. A string is *square-free* if none of its substring is a square. Given a string S ,
 129 a *contraction* is the reverse of a tandem duplication. That is, it takes a square string
 130 XX contained in S and deletes one of the two copies of X . We write $T \mapsto S$ if there
 131 exist strings A, B, X such that $T = AXXB$ and $S = AXB$. We also define $T \mapsto_k S$
 132 and $T \mapsto_* S$ for contractions analogously as for TDs. Observe that by the symmetry
 133 of duplications and contractions, $T \mapsto_k S$ if and only if $S \Rightarrow_k T$ and $T \mapsto_* S$ if and

¹In other words, if we were to state the maximization version of the Cost-Effective Subgraph problem, p would be the value to maximize. The minimization version, however, is more convenient to use for our needs.

134 only if $S \Rightarrow_* T$. When there is no possible confusion, we will sometimes write $T \mapsto S$
 135 instead of $T \mapsto_* S$.

136 We have the following problem.

137 The Tandem Duplication (TD) problem:

138 **Input:** Two strings S and T over the same alphabet Σ and an integer k .

139 **Question:** Is $\text{dist}_{TD}(S, T) \leq k$?

140 In the Exemplar-TD variant of this problem, S is required to be exemplar. In
 141 either variant, we may call S the *source string* and T the *target string*. We will often
 142 use the fact that S and T form a YES instance if and only if T can be transformed
 143 into S by a sequence of at most k contractions. See Fig.1 for a simple example.

	<i>Sequence</i>	<i>Operations</i>
<i>Sequence</i> $T =$	$\langle a, c, \underline{g}, g, a, c, g \rangle$	<i>contraction on</i> $\langle g, g \rangle$
	$\langle a, c, \underline{g}, a, c, g \rangle$	<i>contraction on</i> $\langle a, c, g, a, c, g \rangle$
<i>Sequence</i> $S =$	$\langle a, c, g \rangle$	

FIG. 1. An example for transforming sequence T to S by two contractions. The corresponding sequence of TDs from S to T would duplicate a, c, g , and then duplicate the first g .

144 We recall that although we study the minimization problem here, it is unknown
 145 whether the question $S \Rightarrow_* T$ can be decided in polynomial time. Nonetheless, our
 146 NP-hardness reduction applies to ‘promise’ instances in which $S \Rightarrow_* T$ always holds.

147 **3. NP-hardness of Exemplar-TD.** To facilitate the presentation of our hard-
 148 ness proof, we first make an intermediate reduction using the Cost-Effective Subgraph
 149 problem, which we will then reduce to the promise version of the Exemplar-TD prob-
 150 lem.

151 **The Cost-Effective Subgraph problem.** Suppose we are given a graph $G =$
 152 (V, E) and an integer cost $c \in \mathbb{N}_{>0}$. For $X \subseteq V$, let $E(X) = \{uv \in E : u, v \in X\}$
 153 denote the edges inside of X . The *cost* of X is defined as

$$\text{cost}(X) = c \cdot (|E(G)| - |E(X)|) + |X| \cdot |E(X)|.$$

154 The Cost-Effective Subgraph problem asks for a subset X of minimum cost. In
 155 the decision version of the problem, we are given an integer r and we want to know
 156 if there is a subset X whose cost is at most r . Observe that $X = \emptyset$ or $X = V$ are
 157 possible solutions.

158 The idea is that each edge “outside” of X costs c and each edge “inside” costs
 159 $|X|$. Therefore, we pay for each edge not included in X , but if X gets too large,
 160 we pay more for edges in X . We must therefore find a balance between the size of
 161 X and its number of edges. The connection with the TD problem can be roughly
 162 described as follows: in our reduction, we will have many substrings which need to
 163 be deleted through contractions. We will have to choose an initial set of contractions
 164 X and then, each substring will have two ways to be contracted: one way requires c
 165 contractions, and the other requires $|X|$.

166 An obvious solution for a Cost-Effective Subgraph is to take $X = \emptyset$, which is of
 167 cost $c|E(G)|$. Another formulation of the problem could be whether there is a subset
 168 X of cost at most $c|E(G)| - p$, where p can be seen as a “profit” to maximize. Treating
 169 c and p as parameters, we show the NP-hardness and W[1]-hardness in parameters
 170 $c + p$ of the Cost-Effective Subgraph problem (we do not study the parameter r). Our
 171 reduction to the TD problem does not preserve W[1]-hardness and we only use the
 172 NP-hardness in this paper, but the W[1]-hardness might be of independent interest.

173 Before proceeding, we briefly argue the relevance of parameter c in the W[1]-
 174 hardness. If c is a fixed constant, then we may assume that any solution X satisfies
 175 $|X| \leq c$. This is because if $|X| > c$, every edge included in X will cost more than c
 176 and putting $X = \emptyset$ yields a lower cost. Thus for fixed c , it suffices to brute-force every
 177 subset X of size at most c and we get a $n^{O(c)}$ time algorithm. Our W[1]-hardness
 178 shows that it is difficult to remove this exponential dependence between n and c .

179 **THEOREM 3.1.** *The Cost-Effective Subgraph problem is NP-hard and W[1]-hard*
 180 *for parameter $c + p$.*

181 *Proof.* We reduce from CLIQUE. In this classic problem, we are given a graph
 182 G and an integer k , and must decide whether G contains a clique of size at least k ,
 183 where a clique is a set of vertices in which every pair shares an edge. This problem is
 184 NP-hard [19] and also W[1]-hard in parameter k [10]. We will assume that k is even
 185 (which does not alter either hardness results).

186 Let (G, k) be a CLIQUE instance, letting $n := |V(G)|$ and $m := |E(G)|$. The
 187 graph in our Cost-Effective Subgraph instance is also G . We set the cost $c = 3k/2$,
 188 which is an integer since k is even, and set

$$189 \quad r := c \left(m - \binom{k}{2} \right) + k \binom{k}{2} = cm + \binom{k}{2} (k - c) = cm - \frac{k}{2} \binom{k}{2}.$$

191 We ask whether G admits a subgraph X satisfying $cost(X) \leq r$. We show that
 192 (G, k) is a YES instance to CLIQUE if and only if G contains a set $X \subseteq V(G)$ of cost
 193 at most r . This will prove both NP-hardness and W[1]-hardness in $c + p$ (noting that
 194 here $p = k/2 \binom{k}{2}$).

195 The forward direction is easy to see. If G is a YES instance, it has a clique X of
 196 size exactly k . Since $|E(X)| = \binom{k}{2}$, the cost of X is precisely r .

197 Let us consider the converse direction. Assume that (G, k) is a NO instance of
 198 CLIQUE. Let $X \subseteq V(G)$ be any subset of vertices. We will show that $cost(X) > r$.
 199 There are 3 cases to consider depending on $|X|$.

200 *Case 1:* $|X| = k$. Since G is a NO instance, X is not a clique and thus $|E(X)| =$
 201 $\binom{k}{2} - h$, where $h > 0$. We have that $cost(X) = c(m - \binom{k}{2} + h) + k(\binom{k}{2} - h) =$
 202 $cm + \binom{k}{2}(k - c) + h(c - k) = r + h(c - k)$. Since $c > k$ and $h > 0$, the cost of X is
 203 strictly greater than r .

204 *Case 2:* $|X| = k + l$ for some $l > 0$. Denote $|E(X)| = \binom{k+l}{2} - h$, where $0 \leq h \leq \binom{k+l}{2}$
 205 The cost of X is

$$206 \quad cost(X) = c \left(m - \binom{k+l}{2} + h \right) + (k+l) \left(\binom{k+l}{2} - h \right)$$

$$207 \quad = cm + \binom{k+l}{2} (k+l - c) + h(c - k - l)$$

$$208 \quad = cm + \binom{k+l}{2} \left(l - \frac{k}{2} \right) + h \left(\frac{k}{2} - l \right).$$

209

210 Considering the difference

$$\begin{aligned}
211 \quad \text{cost}(X) - r &= \binom{k+l}{2} \left(l - \frac{k}{2} \right) + h \left(\frac{k}{2} - l \right) - \left(-\frac{k}{2} \right) \binom{k}{2} \\
212 &= \frac{3kl^2}{4} - \frac{kl}{4} + \frac{l^3}{2} - \frac{l^2}{2} + h \left(\frac{k}{2} - l \right), \\
213
\end{aligned}$$

214 if $k/2 - l \geq 0$, then the difference is clearly above 0 regardless of h , and then $\text{cost}(X) >$
215 r as desired. Thus we may assume that $k/2 - l < 0$. In this case, we may further
216 assume that $h = \binom{k+l}{2}$, as this minimizes the difference. But in this case,

$$\text{cost}(X) = cm + \binom{k+l}{2} \left(l - \frac{k}{2} \right) + \binom{k+l}{2} \left(\frac{k}{2} - l \right) = cm > r,$$

217 which concludes this case.

218 *Case 3:* $|X| = k - l$, with $l > 0$. If $k = l$, then $X = \emptyset$ and $\text{cost}(X) = cm > r$. So we
219 assume $k > l$. Put $|E(X)| = \binom{k-l}{2} - h$, where $h \geq 0$. We have

$$\begin{aligned}
220 \quad \text{cost}(X) &= c \left(m - \binom{k-l}{2} + h \right) + (k-l) \left(\binom{k-l}{2} - h \right) \\
221 &= cm + \binom{k-l}{2} (k-l-c) + h(c-k+l) \\
222 &= cm + \binom{k-l}{2} \left(-\frac{k}{2} - l \right) + h \left(\frac{k}{2} + l \right). \\
223
\end{aligned}$$

224 The difference with this cost and r is

$$\begin{aligned}
225 \quad \text{cost}(X) - r &= \binom{k-l}{2} \left(-\frac{k}{2} - l \right) + h \left(\frac{k}{2} + l \right) - \left(-\frac{k}{2} \right) \binom{k}{2} \\
226 &= \frac{3kl^2}{4} + \frac{kl}{4} - \frac{l^3}{2} - \frac{l^2}{2} + h \left(\frac{k}{2} + l \right) \\
227 &> \frac{1}{4}(3l^3 + l^2) - \frac{1}{2}(l^3 + l^2) \geq 0, \\
228
\end{aligned}$$

229 the latter since $k > l \geq 1$. Again, it follows that $\text{cost}(X) > r$. \square

Reduction to Exemplar-TD (promise version). Since the reduction is somewhat technical, we provide an overview of the techniques that we will use. Let (G, c, r) be a **Cost-Effective Subgraph** instance where c is the cost and r the optimization value, and with vertices $V(G) = \{v_1, \dots, v_n\}$. Here we present a simplified idea of the reduction, and leave for later the technical modifications required to make this idea work. We will construct strings S and T and argue on the number of contractions to go from T to S . We would like our source string to be $S = x_1x_2 \dots x_n$, where each x_i is a distinct character that corresponds to vertex v_i . Then, the T string will handle the relationship between vertices and edges. First, T begins with a substring S' obtained by doubling every x_i , i.e., $S' = x_1x_1x_2x_2 \dots x_nx_n$. Second, each edge $e_i \in E(G)$ will have a corresponding substring E_i in T . We call E_i the *substring gadget* of e_i , and its exact content will be detailed later. Our goal is to put $T = S'E_1E_2 \dots E_m$. It follows that, to contract T into S , we must get rid of both the doubled x_i 's and get rid of the E_i 's. We make it so that we first want to contract some, but not necessarily all, of the doubled x_i 's of S' , resulting in another string S'' . The choice of x_i 's to contract in this

first step corresponds to vertices in a Cost-Effective Subgraph. Let t be the number of x_i 's contracted from S' to S'' . For instance, we could have $S'' = x_1x_1x_2x_3x_3x_4x_5x_5$, where only x_2 and x_4 were contracted, and thus $t = 2$. After T is transformed to $S''E_1 \dots E_m$, we then force each E_i to use S'' to contract it. For $m = 3$, a contraction sequence that we would like to enforce would take the form

$$\underline{S'E_1E_2E_3} \rightsquigarrow \underline{S''E_1E_2E_3} \rightsquigarrow \underline{S''E_2E_3} \rightsquigarrow \underline{S''E_3} \rightsquigarrow \underline{S''} \rightsquigarrow S,$$

230 where we underline the substring affected by contractions at each step. We make it
 231 so that when contracting $S''E_iE_{i+1} \dots E_m$ into $S''E_{i+1} \dots E_m$, we have two options.
 232 Suppose that v_j, v_k are the endpoints of edge e_i . If, in S'' , we had chosen to contract x_j
 233 and x_k , we can contract E_i using a sequence of t moves. Otherwise, we must contract
 234 E_i using another more costly sequence of c moves. The total cost to eliminate the E_i
 235 gadgets will be $c(m - e) + te$, where e is the number of edges that can be contracted
 236 using the first choice, i.e. for which both endpoints were chosen in S'' .

237 Unfortunately, constructing S' and the E_i 's to implement the above idea is not
 238 straightforward. The main difficulty lies in forcing an optimal solution to behave as
 239 we describe, i.e. enforcing going from S' to S'' first, enforcing the E_i 's to use S'' ,
 240 and enforcing the two options to contract E_i with the desired costs. In particular,
 241 we must replace the x_i 's by carefully constructed substrings X_i . We must also insert
 242 in T many consecutive copies of the $E_1 \dots E_m$ substring, since this will multiply the
 243 impact of choosing a suboptimal set of contractions of S' in the first step. This in turn
 244 allows us to argue that, in an optimal sequence of contractions, "enough" $E_1 \dots E_m$
 245 copies behave as we want in the above form, which is sufficient for our purposes.

246 Another important remark on the construction of the E_i gadget is that they
 247 should not be able to contract each other cheaply. That is, we want only the S'' from
 248 the above form to be able to remove them, unless we are willing to pay a high cost to
 249 avoid it. In order to achieve this, each E_i will have two parts $B_i \dots \Delta B_{2p} \dots \Delta$, each
 250 ending with a separation symbol Δ . These are designed so that each B_i is different,
 251 and the first part of E_{i-1} cannot affect the first part of E_i (because they start with
 252 distinct B_{i-1} and B_i , respectively). However, the second part of E_{i-1} will be able to
 253 remove all of E_i , but at a high cost. We can do this by making each E_i start with
 254 $B_iB_{i-1} \dots B_2B_1^*B_0$, and making the second part start with $B_{2p} \dots B_2B_1B_0$, where a
 255 high number of contractions are needed to turn B_1^* into B_1 .

256 We now proceed with the technical details. To facilitate the readers, we will also
 257 present a detailed example after the proof of the following theorem.

258 **THEOREM 3.2.** *The Exemplar-TD problem is NP-complete, even if for the given*
 259 *string S and T , $S \Rightarrow_* T$ is guaranteed to hold.*

260 *Proof.* To see that the problem is in NP, note that $dist_{TD}(S, T) \leq |T|$ since each
 261 contraction from T to S removes at least character. Thus a sequence of contractions
 262 can serve as a certificate, has polynomial size and is easy to verify.

263 For hardness, we reduce from the Cost-Effective Subgraph problem, which has been
 264 shown NP-hard in Theorem 3.1. Let (G, c, r) be an instance of Cost-Effective Subgraph,
 265 letting $n := |V(G)|$ and $m := |E(G)|$. Here c is the "outsider edge" cost and we ask
 266 whether there is a subset $X \subseteq V(G)$ such that $c(m - |E(X)|) + |X||E(X)| \leq r$. We
 267 denote $V(G) = \{v_1, \dots, v_n\}$ and $E(G) = \{e_1, \dots, e_m\}$. The ordering of vertices and
 268 edges is arbitrary but remains fixed for the remainder of the proof. For convenience,
 269 we allow the edge indices to loop through 1 to m , and so we put $e_i = e_{i+lm}$ for any
 270 integer $l \geq 0$. Thus we may sometimes refer to an edge e_h with an index $h > m$,
 271 meaning that e_h is actually the edge $e_{((h-1) \bmod m)+1}$.

272 **The construction.** Now we show how to construct S and T . The reader should
 273 bear in mind that all strings described in the construction can be generated in time
 274 polynomial with respect to n and m , and we will not reiterate this point. First let
 275 $d = m + 1$ and $p = m(n + m)^{10}$. The exact values of d and p are not crucial and will
 276 only refer to them when needed: for the most part, it is enough to think of d and p
 277 as simply “large enough”, but with p much larger than d . It is however important to
 278 note that p is a multiple of m . Roughly speaking, d controls the cost for the number
 279 of contractions chosen in the first step discussed above, and p/m is the number of
 280 $E_1 \dots E_m$ copies in T . For later reference, the value of k that we will use in the
 281 reduction to the TD instance is $k = p/m \cdot d(r + nm) + 4cdn$.

282 Instead of having a single character x_i corresponding to each vertex v_i of G , we
 283 have a substring X_i of length d corresponding to v_i . It is then the characters of
 284 each X_i that we double. More precisely, for each $i \in [n]$, define an exemplar string
 285 $X_i = x_i^1 x_i^2 \dots x_i^d$ of length d . Moreover, create enough characters so that no two X_i
 286 strings contain a character in common. Let $X_i^d = x_i^1 x_i^1 x_i^2 x_i^2 \dots x_i^d x_i^d$ be obtained from
 287 X_i by doubling every character, and notice that $\text{dist}_{TD}(X_i, X_i^d) = d$.

288 Also define the strings

$$289 \quad \mathcal{X} = X_1 X_2 \dots X_n, \quad \mathcal{X}^d = X_1^d X_2^d \dots X_n^d.$$

291 Then for each $j \in \{0, 1, \dots, 2p\}$, define an exemplar string B_j . Ensure that no B_j
 292 contains a character from an X_i string, and no two B_j 's contain a common character.
 293 The B_j strings can consist of a single character, with the exception of B_0 and B_1
 294 which are special. We assume that for B_0 and B_1 , we have strings B_0^* and B_1^* such
 295 that

$$296 \quad \text{dist}_{TD}(B_0, B_0^*) = dc + 2d - 2,$$

$$297 \quad \text{dist}_{TD}(B_1, B_1^*) = dn + 2d - 1.$$

299 The B_0 and B_0^* strings can be constructed in the same manner as the X_i and X_i^d
 300 strings, i.e. using a doubling trick. The same applies to B_1 and B_1^* . The B_j 's are the
 301 building blocks of larger strings. For each $q \in [2p]$, define

$$302 \quad \mathcal{B}_q = B_q B_{q-1} \dots B_2 B_1 B_0, \quad \mathcal{B}_q^0 = B_q B_{q-1} \dots B_2 B_1 B_0^*,$$

$$303 \quad \mathcal{B}_q^1 = B_q B_{q-1} \dots B_2 B_1^* B_0, \quad \mathcal{B}_q^{01} = B_q B_{q-1} \dots B_2 B_1^* B_0^*.$$

305 These strings are used as “blockers” and prevent certain contractions from happening.
 306 Note that \mathcal{B}_q^0 and \mathcal{B}_q^1 can be turned into \mathcal{B}_q using $dc + 2d - 2$ contractions and $dn + 2d - 1$
 307 contractions, respectively. Moreover, \mathcal{B}_q^{01} can be turned into \mathcal{B}_q^0 using $dn + 2d - 1$
 308 contractions and into \mathcal{B}_q^1 using $dc + 2d - 2$ contractions.

309 Finally, define a new additional character Δ , which will be used to separate some
 310 of the components of our string. We can now define S by putting

$$311 \quad S = \mathcal{B}_{2p} \mathcal{X} \Delta = B_{2p} B_{2p-1} \dots B_2 B_1 B_0 X_1 X_2 \dots X_n \Delta.$$

313 It follows from the definitions of \mathcal{B}_{2p} , \mathcal{X} and Δ that S is exemplar.

314 Next, we define T . But first the E_i gadgets must be defined. For each edge
 315 $e_q = v_i v_j$ with $q \in [p]$ whose endpoints are v_i and v_j , define

$$316 \quad \mathcal{X}_{e_q} = X_1^d \dots X_{i-1}^d X_i X_{i+1}^d \dots X_{j-1}^d X_j X_{j+1}^d \dots X_n^d.$$

318 Thus in \mathcal{X}_{e_q} , all X_h substrings are turned into X_h^d , except X_i and X_j . Now for
 319 $i \in [p]$, define

$$E_i := \mathcal{B}_i^{01} \mathcal{X}_{e_i} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta,$$

320 which we will call the *edge gadget*. Define T as

$$\begin{aligned} 321 \quad T &= \mathcal{B}_{2p}^0 \mathcal{X}^d \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta E_1 E_2 \dots E_p \\ 322 \quad &= \mathcal{B}_{2p}^0 \mathcal{X}^d \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta [\mathcal{B}_1^{01} \mathcal{X}_{e_1} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] [\mathcal{B}_2^{01} \mathcal{X}_{e_2} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] \dots [\mathcal{B}_p^{01} \mathcal{X}_{e_p} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta]. \end{aligned}$$

324 We add brackets for clarity only — they indicate the distinct E_i substrings, but the
 325 brackets are not actual characters of T . The idea is that T starts with $S' = \mathcal{B}_{2p}^0 \mathcal{X}^d \Delta$,
 326 a modified S in which \mathcal{B}_{2p} becomes \mathcal{B}_{2p}^0 and the X_i substrings are turned into X_i^d .
 327 This \mathcal{X}^d substring serves as a choice of vertices in our cost-effective subgraph. Each
 328 edge e_i has a “gadget substring” $E_i = \mathcal{B}_i^{01} \mathcal{X}_{e_i} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta$. Since p is a multiple of m ,
 329 the sequence of edge gadgets $E_1 E_2 \dots E_m$ is repeated p/m times. Our goal to go from
 330 T to S is to get rid of all these edge gadgets by contractions. Note for later reference
 331 that each pair E_i and E_{i+1} start with different character, namely those of B_i and of
 332 B_{i+1} , respectively.

333 **The hardness proof.** We now show that G admits a subset of vertices W of cost at
 334 most r if and only if T can be contracted to S using at most $p/m \cdot d(r + nm) + 4cdn$
 335 contraction operations. It can be deduced from the (\Rightarrow) direction that $T \rightsquigarrow_* S$ holds.

336 (\Rightarrow) Suppose that G admits a subset of vertices W of cost at most r . Thus
 337 $c(m - |E(W)|) + |W| \cdot |E(W)| \leq r$. To go from T to S , first consider an edge e_i
 338 that does not have both endpoints in W , with $i \in [p]$. We show how to get rid
 339 of the gadget substring E_i for e_i using $dn + dc$ contractions. Note that T contains
 340 the substring $\mathcal{B}_{2p}^1 \mathcal{X} \Delta E_i = \mathcal{B}_{2p}^1 \mathcal{X} \Delta [\mathcal{B}_i^{01} \mathcal{X}_{e_i} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta]$, where brackets surround the E_i
 341 occurrence that we want to remove (note that here for $i > 1$, the prefix $\mathcal{B}_{2p}^1 \mathcal{X} \Delta$ is the
 342 suffix of the previous E_{i-1} gadget, and for $i = 1$, it is the suffix of the starting block
 343 of T). We can first contract the \mathcal{B}_i^{01} of E_i to \mathcal{B}_i^1 using $dc + 2d - 2$ contractions, then
 344 contract \mathcal{X}_{e_i} to \mathcal{X} using $d(n - 2)$ contractions. The result is the $\mathcal{B}_{2p}^1 \mathcal{X} \Delta [\mathcal{B}_i^1 \mathcal{X} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta]$
 345 substring, which becomes $\mathcal{B}_{2p}^1 \mathcal{X} \Delta$ using two contractions (see below). This sums to
 346 $dc + 2d - 2 + d(n - 2) + 2 = dc + dn$ operations. More visually, the sequence of
 347 contractions works as follows (as before, brackets indicate the E_i substring and what
 348 remains of it, and the underlines are there to emphasize the substrings that participate
 349 in the contractions(s)):

$$\begin{aligned} 350 \quad & \mathcal{B}_{2p}^1 \mathcal{X} \Delta \left[\underline{\mathcal{B}_i^{01} \mathcal{X}_{e_i} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta} \right] && (dc + 2d - 2 \text{ contractions}) \\ 351 \quad & \rightsquigarrow \mathcal{B}_{2p}^1 \mathcal{X} \Delta \left[\underline{\mathcal{B}_i^1 \mathcal{X}_{e_i} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta} \right] && (d(n - 2) = dn - 2d \text{ contractions}) \\ 352 \quad & \rightsquigarrow \mathcal{B}_{2p}^1 \mathcal{X} \Delta \left[\underline{\mathcal{B}_i^1 \mathcal{X} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta} \right] \\ 353 \quad & = B_{2p} B_{2p-1} \dots B_{i+1} \underline{\mathcal{B}_i^1 \mathcal{X} \Delta} \left[\underline{\mathcal{B}_i^1 \mathcal{X} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta} \right] && (1 \text{ contraction}) \\ 354 \quad & \rightsquigarrow B_{2p} B_{2p-1} \dots B_{i+1} \underline{\mathcal{B}_i^1 \mathcal{X} \Delta} \left[\underline{\mathcal{B}_{2p}^1 \mathcal{X} \Delta} \right] \\ 355 \quad & = \underline{\mathcal{B}_{2p}^1 \mathcal{X} \Delta} \left[\underline{\mathcal{B}_{2p}^1 \mathcal{X} \Delta} \right] && (1 \text{ contraction}) \\ 356 \quad & \rightsquigarrow \underline{\mathcal{B}_{2p}^1 \mathcal{X} \Delta}. \end{aligned}$$

358 This sequence of $dn + dc$ contractions effectively removes the E_i substring gadget.
 359 Observe that after applying this sequence, it is still true that every remaining E_j

360 gadget substring is preceded by $\mathcal{B}_{2p}^1 \mathcal{X} \Delta$. We may therefore repeatedly apply this
 361 contraction sequence to every e_i not contained in W (including those E_i gadgets for
 362 which $i > m$). This procedure is thus applied to $p/m \cdot (m - |E(W)|)$ gadgets. We
 363 assume that we have done so, and that every e_i for which the E_i gadget substring
 364 remains is in W . Call the resulting string T' . Thus

$$365 \quad T' = \mathcal{B}_{2p}^0 \mathcal{X}^d \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta E_{a_1} \dots E_{a_l}$$

366 for some indices a_1, \dots, a_l corresponding to edges inside W .

367 Now, let \mathcal{X}_W be the substring obtained from \mathcal{X}^d by contracting, for each $v_i \in W$,
 368 the string X_i^d to X_i . We assume that we have contracted the \mathcal{X}^d substring of T' to
 369 \mathcal{X}_W , which uses $d|W|$ contractions (note that there is only one occurrence of \mathcal{X}^d in
 370 T' , namely right before the first Δ). Call

$$371 \quad T' = \mathcal{B}_{2p}^0 \mathcal{X}_W \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta E_{a_1} \dots E_{a_l}$$

372 the resulting string. At this point, for every E_i substring gadget that remains, where
 373 E_i corresponds to edge $e_i = v_j v_k$, \mathcal{X}_W contains the substrings X_j and X_k (instead of
 374 X_j^d and X_k^d).

375 Let i be the smallest integer for which the e_i substring gadget E_i is still in T''
 376 (i.e., $i = a_1$ using the notation for T'' just above). This is the leftmost edge gadget
 377 still in T'' , meaning that T'' has the prefix

$$378 \quad \mathcal{B}_{2p}^0 \mathcal{X}_W \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta [\mathcal{B}_i^{01} \mathcal{X}_{e_i} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] \dots,$$

380 where brackets indicate the E_i substring. To remove E_i , first contract \mathcal{B}_i^{01} to \mathcal{B}_i^0 , and
 381 contract \mathcal{X}_{e_i} to X_W (this is possible since $e_i \subseteq W$). The result is

$$382 \quad \mathcal{B}_{2p}^0 \mathcal{X}_W \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta [\mathcal{B}_i^0 \mathcal{X}_W \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] \dots$$

383 One more contraction gets rid of the substring in brackets. This process requires
 384 $dn + 2d - 1 + d(|W| - 2) + 1 = dn + d|W|$ contractions. To recap, the contraction
 385 sequence for E_i does as follows:

$$\begin{aligned} 386 \quad & \mathcal{B}_{2p}^0 \mathcal{X}_W \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta [\mathcal{B}_i^{01} \mathcal{X}_{e_i} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] && (dn + 2d - 1 \text{ contractions}) \\ 387 \quad & \mapsto \mathcal{B}_{2p}^0 \mathcal{X}_W \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta [\mathcal{B}_i^0 \mathcal{X}_{e_i} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] && (d(|W| - 2) \text{ contractions}) \\ 388 \quad & \mapsto \mathcal{B}_{2p}^0 \mathcal{X}_W \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta [\mathcal{B}_i^0 \mathcal{X}_W \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] && (1 \text{ contraction}) \\ 389 \quad & \mapsto \mathcal{B}_{2p}^0 \mathcal{X}_W \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta. \end{aligned}$$

391 After we repeat this for every E_i , all that remains is the string $\mathcal{B}_{2p}^0 \mathcal{X}_W \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta$. We
 392 contract \mathcal{X}_W to \mathcal{X} using $d(n - |W|)$ contractions (in total, going from \mathcal{X}^d to \mathcal{X} required
 393 dn contractions). Then contract \mathcal{B}_{2p}^0 and \mathcal{B}_{2p}^1 to \mathcal{B}_{2p} using $dc + 2d - 2 + dn + 2d - 1 =$
 394 $d(c + n + 4) - 3$ contractions. At this point, we have reached $\mathcal{B}_{2p} \mathcal{X} \Delta \mathcal{B}_{2p} \mathcal{X} \Delta$. One more
 395 contraction of the second half of the string yields S . The summary of the number of

396 contractions made is

$$\begin{aligned}
397 & \frac{p}{m} \cdot (m - |E(W)|) \cdot (dc + dn) + \frac{p}{m} \cdot |E(W)| \cdot (dn + d|W|) + dn + d(c + n + 4) - 3 \\
398 & \leq \frac{p}{m} \cdot (m - |E(W)|) \cdot (dc + dn) + \frac{p}{m} \cdot |E(W)| \cdot (dn + d|W|) + 4cdn \\
399 & = \frac{p}{m} \cdot d \cdot (c + n)(m - |E(W)|) + \frac{p}{m} \cdot d \cdot (n + |W|)|E(W)| + 4cdn \\
400 & = \frac{p}{m} \cdot d \cdot [c(m - |E(W)|) + |W||E(W)| + nm] + 4cdn \\
401 & \leq \frac{p}{m} \cdot d(r + nm) + 4cdn, \\
402 &
\end{aligned}$$

403 as desired.

404 (\Leftarrow) Suppose that T can be turned into S using α contractions, where $\alpha \leq$
405 $p/m \cdot d(r + nm) + 4cdn$. Let C_1, \dots, C_α be a corresponding sequence of contractions.
406 Here, each C_i contraction is given by a pair of positions ranging over both copies of
407 the contracted substring. The idea is to show that, for some integer t , many of the
408 E_i substrings are removed, but only after t of the X_i^d substrings from \mathcal{X}^d have been
409 contracted to X_i . This set of t X_i 's corresponds to the vertices of a cost-effective
410 subgraph. The main components of the proof are to show that 1) each E_i must be
411 removed; 2) no two E_i 's are completely removed by the same contraction; and 3) most
412 (though perhaps not all) E_i require either $dn + dt$ or $dc + dn - 1$ contractions.

413 Denote $T(l)$ as the string obtained from T after applying the first l contractions
414 C_1, \dots, C_l in the sequence, with $T(0) = T$ and $T(\alpha) = S$. A *block* of $T(l)$ is a
415 substring P of $T(l)$ whose last character is Δ , that has only one occurrence of Δ and
416 that is a maximal string with this property (hence in $T(l)$, the first character of P is
417 either preceded by Δ or is the start of $T(l)$). For instance, each E_i substring of T
418 is made of 2 blocks.

419 We need to distinguish between individual occurrences of each character in our
420 strings. A *marker* is a specific occurrence of a character. That is, for string T , each
421 character $T[i]$ is a distinct marker, even if other positions contain the same character.
422 One may think of a marker as having a unique identifier in a string. For concreteness,
423 let us assume that in T , each marker has, as its unique identifier, its position in T ,
424 so that each character is distinguishable. When contracting a substring DD of $T(l)$
425 to the substring D of $T(l + 1)$, we assume that the markers of the second half are
426 deleted. That is, if $T(l) = LDDR$ and $T(l + 1) = LDR$, only the markers from
427 the first, leftmost D substring remain. This is without loss of generality, since either
428 choice would lead to the S string. It follows that when going from $T(l)$ to $T(l + 1)$,
429 some markers might change position but they remain “the same” marker (i.e. their
430 identifier remains the same). This lets us assert that each marker of $T(l)$ is a marker
431 that was originally in T .

432 For a substring P of T , we say that P is *removed in $T(l)$* if none of the markers
433 of P are in $T(l)$. We say that P is removed if there exists some $T(l)$, with $1 \leq l \leq \alpha$,
434 in which P is removed.

435 CLAIM 1. *Each E_i substring must be removed in $T(\alpha)$.*

436 *Proof.* Consider the first, leftmost block $\mathcal{B}_{2p}^0 \mathcal{X}^d \Delta$ of T . Notice that the Δ marker
437 of this block is the first occurrence of Δ in T . Also note that for any $T(l)$, the Δ
438 marker from this block must also be present in $T(l)$, since there is no way to remove it,
439 by our way of always deleting the rightmost copy in contractions. Now, if we assume
440 that some E_i is not removed from $T(\alpha)$, some marker x of E_i occurs in $T(\alpha) = S$.

441 Since x occurs to the right of the first Δ in T , x also occurs to the right of the first Δ
 442 in S , because contractions cannot change the relative ordering of two markers. This
 443 is a contradiction, since Δ is the last character of S . \square

444 For what follows, notice that in T , $E_i = \mathcal{B}_i^{01} \mathcal{X}_{e_i} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta$ has two blocks. We write
 445 $E'_i = \mathcal{B}_i^{01} \mathcal{X}_{e_i} \Delta$ to denote the first block of E_i . For the rest of the proof, we will mainly
 446 focus on how the sequence of contractions removes the first block E'_i . This gives us
 447 sufficient bounds on the number of contractions required to remove each E_i , which
 448 simplifies the analysis. We let $E'_i(l)$ be the substring of $T(l)$ formed by all the markers
 449 that belong to E'_i , noting that $E'_i(l)$ is possibly the empty string or a subsequence of
 450 E'_i . For $a \in \{0, 1, 01\}$, a block $BX\Delta$ is called a $\mathcal{B}_i^a \mathcal{X} \Delta$ -block if $BX\Delta$ is a subsequence
 451 of $\mathcal{B}_i^a \mathcal{X}^d \Delta$ and $\Sigma(BX\Delta) = \Sigma(\mathcal{B}_i^a \mathcal{X}^d \Delta)$. In other words, $BX\Delta$ has all the characters
 452 that occur in $\mathcal{B}_i^a \mathcal{X}^d \Delta$ in the same order, although the number of occurrences of a
 453 character might differ. A string obtained by concatenating an arbitrary number of
 454 $\mathcal{B}_{2p}^1 \mathcal{X} \Delta$ -blocks is called a $\mathcal{B}_{2p}^1 \mathcal{X} \Delta$ -cluster, which could possibly be the empty string.
 455 Using notation borrowed from regular languages, we write $(\mathcal{B}_{2p}^1 \mathcal{X} \Delta)^*$ to denote a
 456 $\mathcal{B}_{2p}^1 \mathcal{X} \Delta$ -cluster.

457 The next claim shows that a specific block structure is preserved during the
 458 contraction process. An important feature of this structure is that each E'_i starts
 459 as a $\mathcal{B}_j^{01} \mathcal{X} \Delta$ -block in T , and remains so until it is completely removed. This means
 460 that until E'_i is removed, all the characters that appear in E'_i remain present (though
 461 perhaps in different numbers of copies).

462 CLAIM 2. For any l , $T(l)$ has the form

$$BX\Delta(\mathcal{B}_{2p}^1 \mathcal{X} \Delta)^* E'_{i_1}(l) (\mathcal{B}_{2p}^1 \mathcal{X} \Delta)^* E'_{i_2}(l) (\mathcal{B}_{2p}^1 \mathcal{X} \Delta)^* \dots E'_{i_h}(l) (\mathcal{B}_{2p}^1 \mathcal{X} \Delta)^*,$$

463 where

- 464 • $BX\Delta$ is a $\mathcal{B}_{2p}^0 \mathcal{X} \Delta$ -block,
- 465 • $1 \leq i_1 < i_2 < \dots < i_h \leq p$,
- 466 • each $(\mathcal{B}_{2p}^1 \mathcal{X} \Delta)^*$ is a $\mathcal{B}_{2p}^1 \mathcal{X} \Delta$ -cluster, and
- 467 • for each $j \in \{i_1, \dots, i_h\}$, $E'_j(l)$ is a $\mathcal{B}_j^{01} \mathcal{X} \Delta$ -block.

468 *Proof.* Notice that the statement is true for $l = 0$, since T has the required form.
 469 Assume the claim is false and let l be the smallest integer for which $T(l)$ is a counter-
 470 example to the claim. Thus we may assume that $T(l-1)$ has the same form as
 471 in the claim statement. Let D be the string that was contracted from $T(l-1)$ to
 472 $T(l)$, so that $T(l-1)$ contained DD as a substring, and the second D substring gets
 473 removed from $T(l-1)$. We shall refer to the first D as the left D and the second
 474 as the right D . If D does not contain a Δ character, then DD is entirely contained
 475 in a single block. Contracting DD cannot remove all occurrences of a character nor
 476 change their order, and thus the above form must be preserved (every $\mathcal{B}_i^a \mathcal{X} \Delta$ -block
 477 will remain a $\mathcal{B}_i^a \mathcal{X} \Delta$ -block). Assume instead that the last character of D is Δ . Then
 478 $DD = D' \Delta D' \Delta$ for some string D' , and removing the second $D' \Delta$ half only removes
 479 entire blocks of $T(l-1)$. As this block cannot be $BX\Delta$ and since each $E'_{i_j}(l-1)$ is
 480 itself a block, this preserves the form of the claim.

481 Therefore, we may assume that the last character of D is not Δ , but that D has
 482 at least one Δ character. Observe that no character from the $BX\Delta$ -block can get
 483 removed by such a contraction, since the left half of DD is kept. It follows that the
 484 first condition of the claim is preserved after contracting DD . It is easy to see that
 485 the second condition is also preserved. For the other two conditions, we have four

486 cases to consider depending on where the right half of DD , i.e. the removed substring,
487 is located in $T(l-1)$.

488 1. The leftmost character removed belongs to a $E'_j(l-1)$ substring of $T(l-1)$.
489 In this case, because D contains a Δ , the right half of DD must contain the Δ
490 of $E'_j(l-1)$. Let b be the first marker of $E'_j(l-1)$. If no confusion is possible,
491 we treat b as both a character and a marker. Note that b is the first character
492 of \mathcal{B}_j^{01} since $E'_j(l-1)$ is a $\mathcal{B}_j^{01}\mathcal{X}\Delta$ -block, by assumption. Also note that the
493 b marker is preceded by Δ in $T(l-1)$. There are two subcases: either this b
494 is the leftmost removed character or not. In the first case, $D = bD'$ for some
495 D' , which we illustrate as follows (we add brackets around the two copies of
496 D , and underline the removed half):

$$T(l-1) = T'[bD'][\underline{bD'}]T'',$$

497 for some strings T' and T'' . Here the second b is the one from $E'_j(l-1)$. Since
498 it is preceded by Δ in $T(l-1)$, this implies that D' (and thus D) ends with
499 a Δ . But we are assuming that D does not end with Δ . Therefore we know
500 that b is not the leftmost character removed from $T(l-1)$.

501 It follows that the b marker belongs to the left half of DD . This case can be
502 illustrated as follows:

$$T(l-1) = T'[D'bD''][\underline{D'bD''}]T'',$$

503 where $D = D'bD''$ for some strings D' and D'' . Here, the first b is the first
504 character from $E'_j(l-1)$. We can argue that D' is not empty: if D' is empty,
505 then D'' contains Δ because we are assuming that D contains Δ . But the
506 first b is the start of $E'_j(l-1)$, and thus the first Δ encountered after this
507 b is the end of the $E'_j(l-1)$ block. In other words, the left $D'bD''$ entirely
508 contains $E'_j(l-1)$, contradicting that the leftmost deleted marker is inside
509 $E'_j(l-1)$.

510 Since D' is not empty, it follows that D must contain Δb as a substring. But
511 there is only one occurrence of Δb in $T(l-1)$, as $E'_j(l-1)$ is the only block
512 that starts with b ($BX\Delta$ and all the blocks in $\mathcal{B}_{2p}^1\mathcal{X}\Delta$ -clusters start with B_{2p} ,
513 and each E'_{i_n} block starts with B_{i_n}). Therefore, $T(l-1)$ cannot contain DD
514 as a substring since it would contain Δb twice, a contradiction.

515 2. The rightmost character removed is in some $E'_j(l-1)$ substring. That is, the
516 last marker of DD is in $E'_j(l-1)$. Since the right D contains Δ , it must
517 contain the Δ that precedes b , the first marker of $E'_j(l-1)$. This means that
518 the D contains Δb as a substring, which has only one occurrence. We get the
519 same contradiction.

520 3. The leftmost and rightmost characters that get removed belong to distinct
521 $(\mathcal{B}_{2p}^1\mathcal{X}\Delta)$ -clusters of $T(l-1)$. This implies the existence of at least one
522 $E'_j(l-1)$ substring in between. The same type of Δb substring argument
523 applies, since the removed D contains the first character of $E'_j(l-1)$ and its
524 preceding Δ .

525 4. The leftmost and rightmost characters that get removed belong to the same
526 $(\mathcal{B}_{2p}^1\mathcal{X}\Delta)$ -cluster. In this case, it is not hard to verify that the result is yet
527 another $(\mathcal{B}_{2p}^1\mathcal{X}\Delta)$ -cluster, which preserves the desired form.

528 The cases above cover every possibility: we have covered the cases where the
529 removed substring begins or ends in a $E'_j(l-1)$, and the two cases where both its
530 extremities end in a cluster. This proves the claim. \square

531 We will say that a contraction C_l affects $E'_i(l)$ if at least one marker of $E'_i(l)$ is
 532 in the substring corresponding to C_l . Recall that C_l spans over both copies of the
 533 contracted substring, and so $E'_i(l)$ could be affected by C_l even if none of its characters
 534 gets removed.

535 CLAIM 3. For any $l \in [\alpha]$, the contraction C_l from $T(l)$ to $T(l+1)$ does not affect
 536 two distinct $E'_i(l)$ and $E'_j(l)$ substrings of $T(l)$.

537 *Proof.* Suppose the claim is false, and let $T(l)[a_1..a_2]$ be the substring of $T(l)$
 538 affected by the contraction, where $T(l)[a_1..a_2] = DD$ for some string D . Assume that
 539 $T(l)[a_1..a_2]$ contains characters from both $E'_i(l)$ and $E'_j(l)$, where $i < j$. Let b_i, b_j be
 540 the first characters of $E'_i(l)$ and $E'_j(l)$, respectively, which are the first characters of
 541 \mathcal{B}_i^{01} and \mathcal{B}_j^{01} , respectively, by Claim 2. Then $T(l)[a_1..a_2]$ must contain the substring
 542 Δb_j , since $E'_j(l)$ occurs later than $E'_i(l)$ in $T(l)$. Since Δb_j occurs only once in $T(l)$ as
 543 argued in the previous claim, Δb_j cannot be a substring of D . This is only possible if
 544 D starts with b_j (and consequently ends with Δ). Now, since $E'_i(l)$ does not contain
 545 b_j , $T(l)[a_1..a_2]$ cannot start with a suffix of $E'_i(l)$. Yet some characters of $E'_i(l)$ are in
 546 $T(l)[a_1..a_2]$, implying that $E'_i(l)$ is entirely in $T(l)[a_1..a_2]$, and furthermore that the
 547 substring Δb_i is in $T(l)[a_1..a_2]$. Again, this substring occurs only once in $T(l)$, and
 548 thus D must start with b_i . But this is impossible since now, D start with b_i and by
 549 b_j , and $b_i \neq b_j$. \square

550 Notice that T has one occurrence of the $\mathcal{X}^d = X_1^d \dots X_n^d$ substring. We will
 551 therefore refer to the \mathcal{X}^d substring of T without ambiguity. For $i \in [n]$, we let $X_i(l)$
 552 denote the substring of $T(l)$ formed by all the markers that belong to the X_i^d substring
 553 of \mathcal{X}^d . We will say that X_i is *activated* in $T(l)$ if $X_i(l) = X_i$. Intuitively speaking,
 554 X_i is activated in $T(l)$ if it has undergone d contractions to turn it from X_i^d into X_i .

555 CLAIM 4. Let $i \in [p]$, and suppose that E'_i is not removed in $T(l-1)$ but is
 556 removed in $T(l)$. Let t be the number of X_i 's that were activated in $T(l-1)$. Suppose
 557 that v_{i_1} and v_{i_2} are the two endpoints of edge e_i .

558 Then the number of contractions that have affected E'_i is at least $dc + dn - 1$ if
 559 X_{i_1} or X_{i_2} is not activated in $T(l-1)$, or at least $\min\{dt + dn, dc + dn - 1\}$ if X_{i_1}
 560 and X_{i_2} are both activated in $T(l-1)$.

561 *Proof.* By Claim 2, in $T(l-1)$, $E'_i(l-1)$ belongs to a $\mathcal{B}_i^{01}\mathcal{X}\Delta$ -block. As $E'_i(l-1)$
 562 gets removed completely after the l -th contraction of some substring DD , it follows
 563 that D must contain a substring that is equal to $E'_i(l-1)$. The right D of the DD
 564 string certainly contains the $E'_i(l-1)$ substring that gets removed, but consider the
 565 copy of $E'_i(l-1)$ in the first D of the DD square. That is, we can represent the
 566 contraction as

$$T'[D_1\hat{E}'_i(l-1)D_2][\underline{D_1E'_i(l-1)D_2}]T'',$$

567 where $D = D_1E'_i(l-1)D_2$ and $\hat{E}'_i(l-1)$ is a substring equal to $E'_i(l-1)$. Since $E'_i(l-1)$
 568 is itself a block, this $\hat{E}'_i(l-1)$ copy is a substring of another block (though $\hat{E}'_i(l-1)$
 569 might not be a block by itself). By Claim 3, the contraction does not affect another
 570 $E'_j(l-1)$, and so there are only two possible blocks that contain $\hat{E}'_i(l-1)$: either it
 571 is $BX\Delta$, which is the $\mathcal{B}_{2p}^0\mathcal{X}\Delta$ -block at the start of $T(l-1)$, or it is a $\mathcal{B}_{2p}^1\mathcal{X}\Delta$ -block
 572 from a cluster preceding $E'_i(l-1)$. We analyze these two cases.

573 Suppose that $\hat{E}'_i(l-1)$ is located in the first block $BX\Delta$ of $T(l-1)$. Note that
 574 since \mathcal{X}_{e_i} contains X_{i_1} and X_{i_2} (as opposed to $X_{i_1}^d$ or $X_{i_2}^d$), X_{i_1} and X_{i_2} must be
 575 activated in $T(l-1)$ for the DD contraction to be possible. Moreover for $E'_i(l-1)$ to

576 be equal to a substring of $BX\Delta$, every other X_j with $j \neq i_1, i_2$ that is activated must
 577 be contracted in $E'_i(l-1)$ (i.e. E'_i contains X_j^d , but must contain X_j in $E'_i(l-1)$).
 578 This requires at least $d(t-2)$ contractions. Moreover, B contains the B_1 substring,
 579 whereas E'_i contains B_1^* . There must have been at least $dn + 2d - 1$ contractions
 580 affecting the \mathcal{B}_i^{01} substring of E'_i . We note that the character doubling strategy forces
 581 all contractions enumerated so far to be distinct. Counting the contraction removing
 582 $E'_i(l-1)$, this implies the existence of $d(t-2) + dn + 2d - 1 + 1 = dn + dt$ contractions
 583 affecting E'_i .

584 If instead $\hat{E}'_i(l-1)$ was located in a $\mathcal{B}_{2p}^1 \mathcal{X}\Delta$ -block, call this block P , then it suffices
 585 to note that P contains B_0 as a substring whereas E'_i contains B_0^* note that B_0 cannot
 586 be modified since it contains only distinct characters, forcing us to transform B_0^* into
 587 B_0). Counting the contraction that removes $E'_i(l-1)$, it follows that at least $dc + 2d - 1$
 588 contractions must have affected E'_i . \square

589 The above shows that there are two types of contractions that can remove E'_i
 590 from $T(l)$. Either it uses the $BX\Delta$ block at the start of $T(l-1)$, or it uses a block
 591 from a preceding $\mathcal{B}_{2p}^1 \mathcal{X}\Delta$ -cluster. We will call the E'_i 's that get removed in the first
 592 manner Type 1, and those that get removed in the second manner Type 2.

593 We would like to show that every Type 1 E'_j gets removed with the same set of
 594 activated X_i 's, but it might not be the case. Rather, our next goal is to show that
 595 “many” E'_j 's of Type 1 use the same activated X_i 's. This is precisely why we created
 596 p/m copies of $E_1 \dots E_m$. For $k \in [p]$, denote by $act(E'_k)$ the set of activated X_i 's
 597 when E'_k gets removed (i.e. when E'_k is not removed from $T(l-1)$ but is removed
 598 from $T(l)$). Recall that the E_i indices range from 1 to p . Let us partition $[p]$ into
 599 intervals of integers $P_a = [am + 1 .. am + m]$, where $a \in \{0, \dots, p/m - 1\}$. We say that
 600 interval P_a is *homogeneous* if, for each $i, j \in P_a$ such that E'_i and E'_j are of Type 1,
 601 $act(E'_i) = act(E'_j)$. In other words, P_a is homogeneous if all the Type 1 E'_i substrings
 602 corresponding to those in P_a are removed with the same set of activated X_i 's.

603 **CLAIM 5.** *There are at least $p/m - 2n$ homogeneous intervals.*

604 *Proof.* Observe that once an X_i is activated, it remains so for the rest of the
 605 contraction sequence. Since there are n of the X_i 's, there are only $n + 1$ possible
 606 values for $act(E'_k)$ (counting the case when none of them are activated). There are
 607 p/m intervals, and it follows that at most $n + 1 \leq 2n$ of them are not homogeneous. \square

608 We can now go on with the final elements of the proof. Define $cost(E'_i)$ as the
 609 number of contractions that affect E'_i . Let P_{a_1}, \dots, P_{a_h} be the set of homogeneous
 610 intervals, $h \geq p/m - 2n$. Choose the P_a interval among those whose sum of corre-
 611 sponding E'_i costs is minimized — in other words choose P_a such that

$$\sum_{i \in P_a} cost(E'_i) = \min_{j \in [h]} \sum_{i \in P_{a_j}} cost(E'_i).$$

612 By Claim 3, no two E'_i 's share their cost, and by the minimality of P_a the total
 613 number of contractions is at least

$$\left(\frac{p}{m} - 2n\right) \sum_{i \in P_a} cost(E'_i)$$

614 We will only bother with these contractions and we make no assumption on the
 615 non-homogeneous intervals. Assume that there is at least one $i \in P_a$ such that E'_i
 616 is of Type 1. Then by Claim 4, $cost(E'_i)$ is either at least $\min\{dc + dn - 1, dt + dn\}$

617 where $t = |\text{act}(E'_i)|$, or $\text{cost}(E'_i)$ is at least $dc + dn - 1$. If $dt + dn \geq dc + dn - 1$, we
 618 may assume that E'_i is of Type 2 since removing E'_i using Type 2 contractions will
 619 not increase its cost. We will therefore assume that if there is at least one E'_i of Type
 620 1 in P_a , then $dt + dn < dc + dn - 1$ and thus $\text{cost}(E'_i) \geq dt + dn$.

621 Now, choose any i in P_a such that E'_i is of Type 1, and let W be the set of vertices
 622 of G corresponding to those in $\text{act}(E'_i)$. That is, $v_j \in W$ if and only if X_j is activated
 623 when E'_i gets removed. If there does not exist an E'_i of Type 1 to choose, then define
 624 $W = \emptyset$. Denote $t := |W|$ and $s := |E(W)|$. We claim that W is a subgraph of G
 625 satisfying $c(m - s) + ts \leq r$.

626 Assume $c(m - s) + ts > r$ (otherwise, we are done). As we are dealing with integers,
 627 this means $c(m - s) + ts \geq r + 1$. We will derive a contradiction on the assumed number
 628 of contractions. For any E'_i where $i \in P_a$, by Claim 4, either an endpoint of e_i is not
 629 in W and $\text{cost}(E'_i) \geq dc + dn - 1$, or e_i is in W and $\text{cost}(E'_i) \geq dt + dn$. Note that we
 630 needed to choose P_a to be homogeneous to guarantee that every Type 1 E'_i uses the
 631 same value of t in the cost $dt + dn$. It follows that the total number of contractions
 632 is at least

$$\begin{aligned}
 633 \quad & \left(\frac{p}{m} - 2n\right) \sum_{i \in P_a} \text{cost}(E'_i) \\
 634 \quad & \geq \left(\frac{p}{m} - 2n\right) [(m - s)(dc + dn - 1) + s(dt + dn)] \\
 635 \quad & = \left(\frac{p}{m} - 2n\right) [d((c + n)(m - s) + s(t + n)) - m + s] \\
 636 \quad & = \left(\frac{p}{m} - 2n\right) [d(c(m - s) + st + nm - ns + ns) - m + s] \\
 637 \quad & = \left(\frac{p}{m} - 2n\right) \cdot d \cdot [c(m - s) + st + nm] + \left(\frac{p}{m} - 2n\right) (s - m) \\
 638 \quad & \geq \left(\frac{p}{m} - 2n\right) \cdot d \cdot [r + 1 + nm] + \left(\frac{p}{m} - 2n\right) (s - m) \\
 639 \quad & = \left(\frac{p}{m} - 2n\right) \cdot d \cdot [r + nm] + \left(\frac{p}{m} - 2n\right) (d + s - m) \\
 640 \quad & = \frac{p}{m} \cdot d(r + nm) - 2dn(r + nm) + \left(\frac{p}{m} - 2n\right) (d + s - m) \\
 641 \quad &
 \end{aligned}$$

642 Now if d and p are large enough, the above is strictly greater $p/m \cdot d(r + nm) + 4cdn$,
 643 leading to a contradiction. Our chosen values $d = m + 1$ and $p = (n + m)^{10}$ easily
 644 verify this. We have therefore shown that W has the desired cost. This concludes the
 645 proof. \square

646 **An example.** In Fig.2, we show a graph G with $W = \{v_4, v_5, v_6\}$. The se-
 647 quence T containing E_1 and E_2 is as follows. We set $d = 4$ (certainly it can
 648 be larger, e.g., 9), $\mathcal{X} = X_1 X_2 X_3 X_4 X_5 X_6$, where $X_i = x_{i1} x_{i2} x_{i3} x_{i4}$. And we set
 649 $X_i^4 = x_{i1} x_{i1} x_{i2} x_{i2} x_{i3} x_{i3} x_{i4} x_{i4}$ and $\mathcal{X}^4 = X_1^4 X_2^4 X_3^4 X_4^4 X_5^4 X_6^4$. For $e_i = (v_j, v_k)$,
 650 $\mathcal{X}_{e_i} = X_1^4 \cdots X_{j-1}^4 X_j X_{j+1}^4 \cdots X_{k-1}^4 X_k X_{k+1}^4 \cdots X_n^4$.

667 Similarly, as $e_2 \in E(W)$, E_2 is contracted as follows.

$$\begin{aligned}
668 & \mathcal{B}_{2p}^0 X_1^4 X_2^4 X_3^4 X_4^4 X_5^4 X_6^4 \Delta B_{2p} B_{2p-1} \cdots B_2 B_1^* B_0 \mathcal{X} \Delta [B_2 B_1^* B_0^* X_1^4 X_2^4 X_3^4 X_4 X_5 X_6^4 \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] \\
669 & = \mathcal{B}_{2p}^0 \mathcal{X}^4 \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta [\mathcal{B}_2^{01} \mathcal{X}_{e_2} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] \\
670 & \rightarrow \mathcal{B}_{2p}^0 X_1^4 X_2^4 X_3^4 X_4 X_5 X_6 \Delta B_{2p} B_{2p-1} \cdots B_2 B_1^* B_0 \mathcal{X} \Delta [B_2 B_1^* B_0^* X_1^4 X_2^4 X_3^4 X_4 X_5 X_6^4 \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] \\
671 & = \mathcal{B}_{2p}^0 \mathcal{X}_W \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta [\mathcal{B}_2^{01} \mathcal{X}_{e_2} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] \\
672 & \rightarrow \mathcal{B}_{2p}^0 X_1^4 X_2^4 X_3^4 X_4 X_5 X_6 \Delta B_{2p} B_{2p-1} \cdots B_2 B_1^* B_0 \mathcal{X} \Delta [B_2 B_1 B_0^* X_1^4 X_2^4 X_3^4 X_4 X_5 X_6^4 \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] \\
673 & = \mathcal{B}_{2p}^0 \mathcal{X}_W \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta [\mathcal{B}_2^0 \mathcal{X}_{e_2} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] \\
674 & \rightarrow \mathcal{B}_{2p}^0 X_1^4 X_2^4 X_3^4 X_4 X_5 X_6 \Delta B_{2p} B_{2p-1} \cdots B_2 B_1^* B_0 \mathcal{X} \Delta [B_2 B_1 B_0^* X_1^4 X_2^4 X_3^4 X_4 X_5 X_6 \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] \\
675 & = \mathcal{B}_{2p}^0 \mathcal{X}_W \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta [\mathcal{B}_2^0 \mathcal{X}_W \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] \\
676 & \rightarrow B_{2p} B_{2p-1} \cdots B_2 B_1 B_0^* X_1^4 X_2^4 X_3^4 X_4 X_5 X_6 \Delta \cdot B_{2p} B_{2p-1} \cdots B_2 B_1^* B_0 X_1 X_2 X_3 X_4 X_5 X_6 \Delta \\
677 & = \mathcal{B}_{2p}^0 \mathcal{X}_W \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta \\
678 & \rightarrow \cdots \text{contract other } E_i \text{ in } E(W) \\
679 & \rightarrow B_{2p} B_{2p-1} \cdots B_2 B_1 B_0^* X_1 X_2 X_3 X_4 X_5 X_6 \Delta \cdot B_{2p} B_{2p-1} \cdots B_2 B_1^* B_0 X_1 X_2 X_3 X_4 X_5 X_6 \Delta \\
680 & = \mathcal{B}_{2p}^0 \mathcal{X} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta \\
682 & = S.
\end{aligned}$$

683 Note that in the reverse direction, we do not necessarily know the optimal way
684 to contract E_k (e.g., when E_k is not in a homogeneous interval). For instance, the
685 optimal solution could contract E_1 as follows

$$\begin{aligned}
686 & \mathcal{B}_{2p}^0 X_1^4 X_2^4 X_3^4 X_4^4 X_5^4 X_6^4 \Delta B_{2p} B_{2p-1} \cdots B_2 B_1^* B_0 \mathcal{X} \Delta [B_1^* B_0^* X_1 X_2^4 X_3 X_4^4 X_5^4 X_6^4 \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] \\
687 & = \mathcal{B}_{2p}^0 \mathcal{X}^4 \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta [\mathcal{B}_1^{01} \mathcal{X}_{e_1} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] \\
688 & \rightarrow \mathcal{B}_{2p}^0 X_1^4 X_2^4 X_3^4 X_4^4 X_5^4 X_6^4 \Delta B_{2p} B_{2p-1} \cdots B_2 B_1^* B_0 \mathcal{X} \Delta [B_1 B_0^* X_1 X_2^4 X_3 X_4^4 X_5^4 X_6^4 \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] \\
689 & = \mathcal{B}_{2p}^0 \mathcal{X}^4 \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta [\mathcal{B}_1^0 \mathcal{X}_{e_1} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] \\
690 & \rightarrow \mathcal{B}_{2p}^0 X_1 X_2^4 X_3 X_4^4 X_5^4 X_6^4 \Delta B_{2p} B_{2p-1} \cdots B_2 B_1^* B_0 \mathcal{X} \Delta [B_1 B_0^* X_1 X_2^4 X_3 X_4^4 X_5^4 X_6^4 \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] \\
691 & = \mathcal{B}_{2p}^0 \mathcal{X}_{e_1} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta [\mathcal{B}_1^0 \mathcal{X}_{e_1} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta] \\
692 & \rightarrow \mathcal{B}_{2p}^0 X_1 X_2^4 X_3 X_4^4 X_5^4 X_6^4 \Delta B_{2p} B_{2p-1} \cdots B_2 B_1^* B_0 X_1 X_2 X_3 X_4 X_5 X_6 \Delta \\
693 & = \mathcal{B}_{2p}^0 \mathcal{X}_{e_1} \Delta \mathcal{B}_{2p}^1 \mathcal{X} \Delta.
\end{aligned}$$

695 Of course, in this way the \mathcal{X}^d after \mathcal{B}_{2p}^0 will lose contents (i.e., some X_i^d will be
696 contracted into X_i). However, the above contractions can only last in at most n
697 homogeneous intervals.

698 In the next section, we show that the hardness result holds even when the alphabet
699 is of size four (fitting that the DNA and sequences are over the set of nucleotides
700 $\{A, C, G, T\}$ and $\{A, C, G, U\}$ respectively).

701 **4. NP-hardness with constant alphabet size.** We reduce the TD problem
702 on arbitrary alphabet to the TD problem on alphabets of size four. Note that this
703 implies that the TD problem is NP-hard on input strings on alphabet of size c for any
704 constant $c \geq 4$.

705 We show that we can encode each individual character $x \in \Sigma$ with a substring
706 X on alphabet of size four and preserve the hardness. The main difficulty is to show

707 that no contraction occurs inside the encoded substrings X , which requires a specific
708 form of square-free strings.

709 LEMMA 4.1. *Let m be an integer. Then there exists an algorithm that takes time*
710 *polynomial in m that outputs pairs of square-free strings $(A_1, B_1), \dots, (A_m, B_m)$ over*
711 *alphabet $\{a, b, c\}$ such that for each $i \in \{1, \dots, m\}$:*

- 712 1. $100i - 3 \leq |A_i| \leq 100i$ and $|B_i| = 1000m - |A_i|$;
- 713 2. $A_i[1] = a$
- 714 3. $B_i[1] = b$ and $B_i[|B_i|] = b$.

715 *Proof.* In [22], Leech describes the following process to generate an infinite ternary
716 square-free string. Define $w_0 = a$, and for $j > 0$, recursively obtain w_j from w_{j-1} by
717 replacing each character with a specified substring of length 13 (the replacement does
718 not matter here, see [22] for details). Each w_j is square-free, and for our purposes we
719 may apply this process and stop when the condition $|w_j| > 2000m$ is met for the first
720 time (note that this takes polynomial time).

721 Assume thus that we generated a square-free string w_j of length above $2000m$.
722 Further assume that the first a occurrence of w_j is at position p , and for $k \geq 1$, define
723 $A(k) = w_j[p..p+k-1]$. Thus each $A(k)$ starts with a and has length k (and we will
724 choose a subset of these to prove the lemma).

725 Similarly, assume that the first b of w_j is at position q . Consider $k \geq 1$. Let
726 $r \geq q+k-1$ be the first position at or after $q+k-1$ such that $w_j[r] = b$. Define
727 $B(k) = w_j[q..r]$. Note that $r \leq q+k+2$, as otherwise $w_j[q+k-1..q+k+2]$ would
728 be a binary string of length 4 and would therefore contain a square. It follows that
729 $k \leq |B(k)| \leq k+3$. Moreover, each $B(k)$ starts and ends with a b .

730 To conclude the lemma, for each $i \in \{1, \dots, m\}$, put $B_i = B(1000m - 100i)$ and
731 $A_i = A(1000m - |B_i|)$. We have $1000m - 100i \leq B_i \leq 1000m - 100i + 3$, and, since
732 $|A_i| = 1000m - |B_i|$, we have $100i - 3 \leq |A_i| \leq 100i$. The other conditions of the
733 lemma follow from our construction. \square

734 We can now proceed to our reduction.

735 THEOREM 4.2. *The TD problem is NP-hard for any $|\Sigma| \geq 4$, where Σ is the*
736 *alphabet of the input strings.*

737 *Proof.* Our reduction is from the TD problem with unbounded alphabet size,
738 which is NP-hard owing to Theorem 3.2. We show that the TD problem is NP-hard
739 when $|\Sigma| = 4$, which clearly implies hardness for any $|\Sigma| \geq 4$ since one can add new
740 dummy character at the end of S and T if $|\Sigma|$ is smaller than the desired value. Let
741 (S, T, k) be an instance of the TD problem, let $\Sigma = \{s_1, \dots, s_m\}$ be the alphabet of
742 S and T , with $m = |\Sigma|$.

743 Let $(A_1, B_1), \dots, (A_m, B_m)$ be pairs of ternary square-free strings that satisfy
744 the conditions of Lemma 4.1, i.e. each $|A_i|$ is between $100i - 3$ and $100i$, $|B_i| =$
745 $1000m - |A_i|$, with A_i starting with a , and B_i starting and ending with b . By the
746 lemma, these can be obtained in polynomial time. Assume that the alphabet of the
747 A_i and B_i strings is $\{a, b, c\}$, and let x be another character.

748 We encode each character $s_i \in \Sigma$ on alphabet $\{a, b, c, x\}$ by defining

$$749 \quad h(s_i) = xcxA_ixB_i$$

750 Note that each $h(s_i)$ string has the same length $1000m + 4$. Also note that each A_i
751 has a distinct length, and thus that each B_i also has a distinct length. This implies
752 that each A_i is matched with a distinct B_i determined by its length (and vice-versa).

753 Moreover, since A_i and B_i are square-free, we can argue that $h(s_i)$ is also square-free.
 754 Indeed, no square can start at the first x of $h(s_i)$ since xc only occurs once in $h(s_i)$;
 755 the square could not be $cxcx$ since the first cx would have to start at the second
 756 character c of $h(s_i)$, and the second cx could only be at the end of A_i , but the latter
 757 has length at least 97, making $cxcx$ impossible; no longer square can start at the
 758 second character c since such a square would need to start with cxa , of which there
 759 is only one occurrence; no square can start at the second x since xa occurs only once;
 760 and no square can be entirely in A_i or entirely in B_i because they are square-free.
 761 Therefore, a square would need to start in A_i and contain the x after it, which is not
 762 possible since there is only one such x . This lets us assert that no tandem duplication
 763 can occur inside an $h(s_i)$.

764 We then encode S and T into S' and T' , respectively, by replacing each character
 765 with their encoding. That is, put

$$766 \quad S' = h(S[1])h(S[2]) \dots h(S[|S|]) \quad T' = h(T[1])h(T[2]) \dots h(T[|T|])$$

767 Observe that in both strings, all occurrences of $xcxa$ mark the start of an encoded
 768 character (since all A_i strings start with a), and occurrences of xb mark a separation
 769 point between some A_i and B_i . In what follows, we say that a string Y *encodes* X if
 770 $Y = h(X[1]) \dots h(X[|X|])$.

771 We show that T can be transformed into S using k contractions if and only if T'
 772 can be turned into S' using k contractions.

773 (\Rightarrow) Assume that T can be transformed into S using k contractions. Let $T =$
 774 $T_0, T_1, \dots, T_k = S$ be the sequence of strings encountered during these contractions.
 775 We define a sequence of contractions from T' to S' in which $T' = T'_0, T'_1, \dots, T'_k = S'$
 776 are the strings encountered during this sequence, such that each T'_i encodes T_i . This
 777 is true when $i = 0$. Assume inductively that T'_{i-1} encodes T_{i-1} , and that the i -th
 778 contraction affects some substring RR of T_{i-1} , i.e. $T_{i-1} = ARRB \rightarrow ARB$. Then
 779 $T'_{i-1} = A'R'R'B'$, where A' encodes A , R' encodes R and B' encodes B . We can
 780 apply $A'R'R'B' \rightarrow A'R'B'$ and obtain T'_i that encodes T_i . It follows that $T'_k = S'$,
 781 as desired.

782 (\Leftarrow) Assume that there is a sequence of k contractions D'_1, \dots, D'_k transforming T'
 783 into S' (assume that the D'_i 's contain the start and end index of the contraction). Let
 784 T'_i be the string obtained after applying the i -th contraction in this sequence, with
 785 $T'_0 = T'$. We construct a sequence of contractions D_1, \dots, D_k transforming T into
 786 S . To achieve this, we construct the D_i contractions by maintaining the following
 787 invariant: by defining T_i as the string obtained after applying D_1, \dots, D_i on T , then
 788 T'_i encodes T_i .

789 This is true for $i = 0$ since $T'_0 = T'$ encodes $T_0 = T$. Now assume inductively
 790 that T'_{i-1} is obtained by encoding each character of T_{i-1} .

791 Assume that D'_i contracts a string RR of T'_{i-1} . Since T'_{i-1} is a concatenation of
 792 $h(s_j)$ strings, the first character of the RR substring must be in some $h(s_j)$ substring
 793 of T'_{i-1} . We consider four possible cases.

794 **Case 1.** The first character of RR is in the xcx substring at the start of $h(s_j) =$
 795 $xcxA_jxB_j$. Thus we may write $xcx = C_1C_2$ such that C_2 is a prefix of RR (note that
 796 C_1 is possibly empty, but not C_2). It is easy to check that R must contain at least a
 797 prefix of A_j . This implies that R contains xa as a substring, i.e. $R = C_2aY$ for some
 798 Y .

799 The contraction $T'_{i-1} \rightarrow T'_i$ has the form

$$800 \quad T'_{i-1} = PC_1C_2aY C_2aYQ \rightarrow PC_1C_2aYQ = T'_i$$

801 for some strings P, Q . Because C_2a contains xa , it must mark the start of an encoding,
 802 since xa can only occur there. Looking at the second copy of C_2a above, we deduce
 803 that Y has C_1 as a suffix. Writing $Y = Y'C_1$, we get

$$804 \quad T'_{i-1} = PC_1C_2aY'C_1C_2aY'C_1Q \rightsquigarrow PC_1C_2aY'C_1Q = T'_i$$

805 We can replace D'_i with the following contraction:

$$806 \quad T'_{i-1} = PC_1C_2aY'C_1C_2aY'C_1Q \rightsquigarrow PC_1C_2aY'C_1Q = T'_i$$

807 Notice that P encodes some prefix \hat{P} of T_{i-1} (P is possibly empty). The C_1C_2aY'
 808 substring that follows also encodes some substring \hat{Y} of T_{i-1} and, since all encodings
 809 have the same length, the second C_1C_2aY' copy that follows encodes the same sub-
 810 string \hat{Y} . This implies that C_1Q , if not empty, starts with xcx and encodes a suffix
 811 \hat{Q} of T_{i-1} . Therefore, $T_{i-1} = \hat{P}\hat{Y}\hat{Y}\hat{Q}$, and it becomes easy to see that if D_i contracts
 812 $\hat{Y}\hat{Y}$, T'_i encodes T_i .

813 **Case 2.** The first character of RR is in the A_j substring of $h(s_j)$. Write $A_j = A_j^1A_j^2$
 814 such that A_j^2 is a prefix of RR . Because A_j is square-free, RR cannot be entirely
 815 contained in A_j . Thus R contains the x after A_j . This x cannot be the last character
 816 of R , as otherwise RR would be fully contained in $h(s_j)$ since $|B_j| > |A_j|$, which is
 817 again not possible since $h(s_j)$ is square-free. Thus the first occurrence of x in R
 818 followed by b , implying that the second R copy must start in the A_l portion of some
 819 other $h(s_l)$ encoding (since this is the only way to have xb as a subtring for the first
 820 x occurrence).

821 Hence, we may assume that $R = A_j^2xB_jxY$ for some non-empty Y . The contrac-
 822 tion has the form

$$823 \quad PxcxA_j^1A_j^2xB_jxY \quad A_j^2xB_jxYQ \rightsquigarrow PxcxA_j^1A_j^2xB_jxYQ$$

824 for some strings P and Q . Since each B_j is associated with a unique A_j , Y must have
 825 $xcxA_j^1$ as a suffix. We may write $Y = Y'xcxA_j^1$ and

$$826 \quad PxcxA_j^1A_j^2xB_jxY'xcxA_j^1 \quad A_j^2xB_jxY'xcxA_j^1Q \rightsquigarrow PxcxA_j^1A_j^2xB_jxY'xcxA_j^1Q$$

827 It becomes easy to see that D'_i can be replaced by the contraction

$$828 \quad PxcxA_j^1A_j^2xB_jxY' \quad xcxA_j^1A_j^2xB_jxY'xcxA_j^1Q \rightsquigarrow PxcxA_j^1A_j^2xB_jxY'xcxA_j^1Q$$

829 and we can apply Case 1.

830 **Case 3.** The first character of RR is the x occurrence in $h(s_j)$ between A_j and
 831 B_j . Thus R starts with xb , implying that the second R copy also starts with the x
 832 occurrence between some A_l and B_l . The contraction $T'_{i-1} \rightsquigarrow T'_i$ has the form

$$833 \quad PxcxA_jxB_jxY \quad xB_jxYQ \rightsquigarrow PxcxA_jxB_jxYQ$$

834 for some strings P, Y, Q . Because B_j is associated only with A_j in the encodings, we
 835 may write $Y = Y'xcxA_j$. We get

$$836 \quad PxcxA_jxB_jxY'xcxA_j \quad xB_jxY'xcxA_jQ \rightsquigarrow PxcxA_jxB_jxY'xcxA_jQ$$

837 We can replace D'_i with an alternate contraction as follows:

$$838 \quad PxcxA_jxB_jxY' \quad xcxA_jxB_jxY'xcxA_jQ \rightsquigarrow PxcxA_jxB_jxY'xcxA_jQ$$

839 and, once again, we can apply Case 1.

840 **Case 4.** The first character of RR is in the B_j substring of $h(s_j)$. Write $B_j = B_j^1 B_j^2$
 841 such that B_j^2 is a prefix of R . We argue that the first R copy cannot be entirely in
 842 B_j^2 . Otherwise, the second R copy would contain characters outside of B_j^2 since B_j is
 843 square-free. That is, the second R copy would contain an x but not the first, which
 844 is impossible.

845 Thus the first R copy contains the x after B_j^2 . We cannot have $R = B_j^2 x$. Indeed,
 846 $B_j^2 x$ is followed by cx , and because $R = B_j^2 x$, the only possible contraction is $RR =$
 847 $B_j^2 xcx = cxcx$. This is not possible, since B_j^2 ends with a b .

848 Therefore, R contains the x after B_j^2 , and also the c after it. The only occurrences
 849 of xc are at the start of encodings, so the contraction has the form

$$850 \quad \underline{PB_j^2 xcY} \quad \underline{B_j^2 xcYQ} \rightsquigarrow \underline{PB_j^2 xcYQ}$$

851 Now, the second R copy starts with a suffix of some B_l , and it happens that B_j and
 852 B_l have the same suffix B_j^2 . We can write $B_l = B_l^1 B_j^2$ such that Y has $x A_l x B_l^1$ as a
 853 suffix. Write $Y = Y' x A_l x B_l^1$ so that

$$854 \quad \underline{PB_j^2 xY' x A_l x B_l^1} \quad \underline{B_j^2 xY' x A_l x B_l^1 Q} \rightsquigarrow \underline{PB_j^2 xY' x A_l x B_l^1 Q} = T'_i$$

855 Recall that A_l is uniquely matched with $B_l = B_l^1 B_j^2$ in the encodings, so Q must have
 856 B_j^2 as a prefix. Write $Q = B_j^2 Q'$ so that

$$857 \quad T'_{i-1} = \underline{PB_j^2 xY' x A_l x B_l^1} \quad \underline{B_j^2 xY' x A_l x B_l^1 B_j^2 Q'} \rightsquigarrow \underline{PB_j^2 xY' x A_l x B_l^1 B_j^2 Q'} = T'_i$$

858 We can replace this contraction with

$$859 \quad T'_{i-1} = \underline{PB_j^2 xY' x A_l x B_l^1 B_j^2} \quad \underline{xcY' x A_l x B_l^1 B_j^2 Q'} \rightsquigarrow \underline{PB_j^2 xY' x A_l x B_l^1 B_j^2 Q'} = T'_i$$

860 We can once again apply Case 1.

861 We have thus shown how to construct a sequence of contractions that transforms
 862 T into a string T_k such that T'_k encodes T_k . Since T'_k encodes S' , it follows that
 863 $T_k = S$, which concludes the proof. \square

864 **5. An FPT algorithm for the exemplar problem.** Let us recall that in the
 865 **Exemplar- k -TD** problem, we receive a string S in which every character is different, a
 866 string T , and an integer k . We want to know whether S can be transformed into T
 867 using at most k tandem duplications. In this section, we show that **Exemplar- k -TD**
 868 can be solved in time $2^{O(k^2)} + poly(n)$ by obtaining a kernel of size $O(k2^k)$, where n
 869 is the length of T .

870 We first note that there is a very simple, brute-force algorithm to solve the k -TD
 871 problem, which is the variant of the TD problem with parameter k , the number of TDs
 872 to turn S into T (including **Exemplar- k -TD** as a particular case). This only establishes
 873 membership in the XP class, but it will be useful to evaluate the complexity of our
 874 kernelization later on.

875 **PROPOSITION 5.1.** *The k -TD problem can be solved in time $O(n^{2k})$, where n is*
 876 *the size of the target string.*

877 *Proof.* Let (S, T) be a given instance of k -TD. Consider the branching algorithm
 878 that, starting from T , tries to contract every substring of the form XX in T and

879 recurses on each resulting string, decrementing k by 1 each time (the branching stops
 880 when S is obtained or when k reaches 0 without attaining S). We obtain a search tree
 881 of depth at most k and degree at most n^2 , and thus it has $O(n^{2k})$ nodes. Visiting the
 882 internal nodes of this search tree only requires enumerating $O(n^2)$ substrings, which
 883 form the set of children of the node. Hence, there is no added computation cost to
 884 consider when visiting a node. \square

885 From now on, we assume that we have an Exemplar- k -TD instance (S, T) , and so
 886 that S is exemplar.

887 Let x and y be two consecutive characters in S (i.e. xy is a substring of S). We
 888 say that xy is an (S, T) -stable pair if in T , every occurrence of x in T is followed by
 889 y and every occurrence of y is preceded by x . That is, the direct successor of every
 890 x character is y , and the direct predecessor of every y character is x . Note that in
 891 this definition, S and T could be any strings. However, for the rest of the section, we
 892 assume that S and T are the input strings.

893 We will show that every (S, T) -stable pair can be replaced by a single character,
 894 and that if T can be obtained from S using at most k tandem duplications, then this
 895 leaves strings of bounded size.

896 We first show that, roughly speaking, stability is maintained by all tandem du-
 897 plications when going from S to T .

898 LEMMA 5.2. *Suppose that $\text{dist}_{TD}(S, T) = k$ and let xy be any (S, T) -stable pair.*
 899 *Let $S = S_0, S_1, \dots, S_k = T$ be any minimum sequence of strings transforming S to T*
 900 *by tandem duplications. Then xy is (S, S_i) -stable for every $i \in [k]$.*

901 *Proof.* Note that xy is trivially (S, S_0) -stable. Assume the lemma is false, and
 902 let S_i be the first of S_1, \dots, S_k that does not satisfy the statement. Then xy is
 903 (S, T) -stable, but xy is not (S, S_i) -stable.

904 We claim that, with our assumption, xy is not (S, S_j) -stable for any $j \in \{i, \dots, k\}$.
 905 As this includes $S_k = T$, this will contradict that xy is (S, T) -stable. We do this
 906 by induction — as a base case, xy is not (S, S_i) -stable so this is true for $j = i$.
 907 Assume that xy is not (S, S_{j-1}) -stable, where $i < j \leq k$. Let D be the duplication
 908 transforming S_{j-1} to S_j (here $D = (a, b)$ contains the start and end positions of the
 909 substring of S_{j-1} to duplicate).

910 Suppose first that xy is not (S, S_{j-1}) -stable because S_{j-1} has an occurrence of x
 911 that is not followed by y . Thus S_{j-1} has an occurrence of xz , with x at some position
 912 p_x , where $z \neq y$. If we assume that xy is (S, S_j) -stable, then a y character must have
 913 appeared after this x from S_{j-1} to S_j . If D contains the xz occurrence, then xz will
 914 still be an occurrence of S_j . If D does not contain the x character, then again xz will
 915 still occur in S_j . It follows that changing the character next to this x is only possible
 916 if the last character duplicated by D is the x at position p_x and the first character
 917 of D is a y . In other words, denoting $S_{j-1} = A_1yA_2xzA_3$ for appropriate A_1, A_2, A_3
 918 substrings, the D duplication must do the following

$$919 \quad A_1yA_2xzA_3 \Rightarrow A_1yA_2xyA_2xzA_3,$$

921 But then, there is still an occurrence of x followed by z , and it follows that xy cannot
 922 be (S, S_j) -stable.

923 So suppose instead that xy is not (S, S_{j-1}) -stable because S_{j-1} has an occurrence
 924 of y preceded by $z \neq x$. Assume the character preceding this y has changed in S_j
 925 and has become x . But one can verify that this is impossible. For completeness, we
 926 present each possible case: either D includes both z and y , includes one of them or

927 none. These cases are represented below, and each one of them leads to an occurrence
 928 of y still preceded by z (the left-hand side represents S_{j-1} and the right-hand side
 929 represents S_j , and the A_i 's are the relevant substrings in each case):

930 Include both: $A_1 \underline{A_2 z y A_3} A_4 \Rightarrow A_1 \underline{A_2 z y A_3 A_2 z y A_3} A_4$,

931 Include z only: $A_1 \underline{A_2 z y A_3} \Rightarrow A_1 \underline{A_2 z A_2 z y A_3}$,

932 Include y only: $A_1 \underline{z y A_2 A_3} \Rightarrow A_1 \underline{z y A_2 y A_2 A_3}$,

933 Include none: $A_1 \underline{A_2 z y A_3} \Rightarrow A_1 \underline{A_2 A_2 z y A_3}$ or $A_1 \underline{z y A_2 A_3} \Rightarrow A_1 \underline{z y A_2 A_2 A_3}$.

935 We have therefore shown that xy cannot be (S, S_j) -stable, and therefore not
 936 (S, T) -stable, which concludes the proof. \square

937 Our next step is to show that no tandem duplication starts or ends between two
 938 characters of a given stable pair xy . This might not be the case in general, but we
 939 show that any duplication that does so can be modified to an equivalent duplication
 940 that does not “cut” any occurrence of xy .

941 Let xy be an (S, T) -stable pair. Let S' be a substring obtained from S by tandem
 942 duplications, and assume that $S'[p_x..p_x + 1] = xy$. Suppose that we then apply a
 943 duplication $D = (c, d)$ on S' , which copies the substring $S'[c..d]$ and pastes it after
 944 position d . Then we say that D cuts xy if one of the following holds:

- 945 • $c = p_x + 1$, in which case we say that D cuts xy to the right;
- 946 • $d = p_x$, in which case we say that D cuts X to the left;

947 In other words, cutting to the right takes the form $AxyBC \Rightarrow AxyByBC$. Cutting
 948 to the left takes the form $ABxyC \Rightarrow ABxByC$. Observe that if D does not cut any
 949 occurrence of xy , and S'' is obtained by applying D on S' , then xy is (S, S'') -stable.

950 LEMMA 5.3. *Suppose that $\text{dist}_{TD}(S, T) = k$ and let xy be an (S, T) -stable pair.
 951 Then there exists a sequence of tandem duplications D_1, \dots, D_k transforming S into
 952 T such that no occurrence of xy gets cut by a D_i .*

953 *In other words, for all $i \in [k]$, the tandem duplication D_i does not cut any occur-
 954 rence of xy in the string obtained by applying D_1, \dots, D_{i-1} to S .*

955 *Proof.* Let D_1^*, \dots, D_k^* be a sequence of tandem duplications transforming S into
 956 T , and for $i \in [k]$, let S_i be the string obtained by applying the first i duplications.
 957 Put $S_0 := S$. We show that any $S_i, i \in [k]$, can be obtained from S_{i-1} by a duplication
 958 D_i that does not cut any occurrence of xy . This proves the lemma, since D_1, \dots, D_k
 959 will form the desired sequence of duplications.

960 Fix $i \in [k]$, and assume that D_i^* cuts at least one xy pair. Notice that D_i^* can cut
 961 at most two occurrences of xy , at most one to the left and at most one to the right.
 962 Also, by Lemma 5.2, we know that every (S, T) -stable xy pair is (S, S_i) -stable. We
 963 have four cases to consider:

- 964 • D_i^* cuts some xy to the right, but no other occurrence to the left. Then we
 965 may write S_{i-1} as $S_{i-1} = AxyBC$ such that D copies the substring yB . This
 966 gives

$$967 \quad S_{i-1} = AxyBC \Rightarrow AxyByBC = S_i.$$

969 But by Lemma 5.2, xy is (S, S_i) -stable. Since S_i has By as a substring, this
 970 must mean that $B = \hat{B}x$ for some substring \hat{B} . Therefore $S_{i-1} = Axy\hat{B}xC$.
 971 Since xy is also (S, S_{i-1}) -stable, this in turn implies that $C = y\hat{C}$ for some
 972 substring \hat{C} , and in fact we get

$$973 \quad S_{i-1} = \underline{Axy\hat{B}xy\hat{C}} \Rightarrow \underline{Axy\hat{B}xy\hat{B}xy\hat{C}} = S_i.$$

975 We can replace D_i^* by a duplication that copies $xy\hat{B}$, i.e.

$$976 \quad S_{i-1} = \underline{Axy\hat{B}xy\hat{C}} \Rightarrow \underline{Axy\hat{B}xy\hat{B}xy\hat{C}} = S_i.$$

978 Recall that x and y are distinct characters because S is exemplar. Since this
 979 duplication starts with xy and copies itself right before another occurrence of
 980 xy , it is clear that it does not cut any occurrence of xy , as desired.

981 • D_i^* cuts some xy to the left, but cuts no string to the right. Then we may
 982 write

$$983 \quad S_{i-1} = \underline{ABxyC} \Rightarrow \underline{ABx\hat{B}xyC} = S_i.$$

985 Similarly as in the previous case, since xy is (S, S_i) -stable, we must have
 986 $B = y\hat{B}$. We are led to deduce that $A = \hat{A}x$. Therefore we have

$$987 \quad S_{i-1} = \underline{\hat{A}xy\hat{B}xyC} \Rightarrow \underline{\hat{A}xy\hat{B}xy\hat{B}xyC} = S_i.$$

989 As before, we could instead duplicate the substring $xy\hat{B}$ occurring right after
 990 \hat{A} .

991 • D_i^* cuts some xy to the left and some xy to the right. We get

$$992 \quad S_{i-1} = \underline{AxyBxyC} \Rightarrow \underline{AxyBxyBxyC} = S_i.$$

994 We could instead duplicate the copy of xyB in S_{i-1} that occurs right after
 995 the A .

996 We have thus shown that if D_i^* cuts some occurrence of one or more of the xy 's,
 997 then D_i^* can be replaced by a duplication D_i that yields the same string S_i as D_i^* . The
 998 only case remaining is when D_i^* does not cut any X_j . In that case, we put $D_i = D_i^*$.
 999 This shows that we can find the claimed sequence D_1, \dots, D_k in which no xy ever
 1000 gets cut. \square

1001 The above implies that we may replace any (S, T) -stable pair xy by a single
 1002 character, since we may assume that characters of xy are always duplicated together
 1003 (assuming, of course, that S is exemplar).

1004 LEMMA 5.4. *Let xy be an (S, T) -stable pair, let z be a character not occurring*
 1005 *in S or T , and let S', T' be obtained by replacing each occurrence of xy in S and T ,*
 1006 *respectively, by z . Then $\text{dist}_{TD}(S, T) = \text{dist}_{TD}(S', T')$.*

1007 *Proof.* Clearly if (S', T') can be solved using at most k duplications, then the
 1008 same applies to (S, T) , since we can modify each duplication that affects z to affect xy
 1009 instead. By Lemma 5.3, the converse also holds: if (S, T) can be solved with at most
 1010 k duplications, we may assume that these duplications never cut an xy occurrence,
 1011 and so these duplications can be applied on (S', T') . \square

1012 It only remains to show that the resulting strings are small enough. This can
 1013 be achieved by repeated application of Lemma 5.4. That is, as long as there is an
 1014 (S, T) -stable pair xy , we can replace xy by a new character z . When this procedure
 1015 terminates, the resulting strings S' and T' have the same TD distance as S and T ,
 1016 but contain no (S, T) -stable pair. We can then rely on the following intuition: a
 1017 tandem duplication can only “break” two stable pairs, so if S' and T' are too long, k
 1018 duplications won’t suffice to break them all.

1019 **LEMMA 5.5.** *Assume that there is no (S, T) -stable pair. If $\text{dist}_{TD}(S, T) \leq k$, then*
 1020 *S has length at most $2k + 1$.*

1021 *Proof.* Let $S = S_0, S_1, \dots, S_k = T$ be any minimum sequence of strings trans-
 1022 forming S to T by tandem duplications. Let n be the length of S . We show by
 1023 induction that, for each $i \in \{0, 1, \dots, k\}$, the number of (S, S_i) -stable pairs is at least
 1024 $n - 2i - 1$. For $i = 0$, there are $n - 1$ (S, S) -stable pairs, so our assertion holds. Now
 1025 assume that there are at least $n - 2(i - 1) - 1 = n - 2i + 1$ (S, S_{i-1}) -stable pairs, with
 1026 $i > 0$. Assume that the duplication transforming S_{i-1} to S_i takes the form

$$1027 \quad S_{i-1} = AxyBuvC \Rightarrow AxyBuyBuvC = S_i$$

1028 for some characters x, y, u, v and substrings A, B, C . The only (S, S_{i-1}) -stable pairs
 1029 that might not be (S, S_i) -stable pairs are xy and uv (since no other substring of
 1030 length two was modified). Note that A, x, v or C might not exist if the duplication
 1031 contains the start or end of S_{i-1} , but in this case the number of stable pairs still
 1032 cannot decrease by more than 2. It follows that the number of (S, S_i) -stable pairs is
 1033 at least $n - 2i + 1 - 2 = n - 2i - 1$, as desired.

1034 By putting $i = k$ and recalling that $S_k = T$, the number of (S, T) -stable pairs is
 1035 at least $n - 2k - 1$. Since this number is 0 by assumption, we get $0 \geq n - 2k - 1$,
 1036 implying $n \leq 2k + 1$. \square

1037 We have thus shown how to transform an instance (S, T) of Exemplar- k -TD to a
 1038 kernel, an equivalent instance (S', T') of size depending only on k .

1039 **THEOREM 5.6.** *An instance (S, T) of Exemplar- k -TD admits a kernel (S', T') in*
 1040 *which $|S'| \leq 2k + 1$ and $|T'| \leq (2k + 1)2^k$.*

1041 *Proof.* As shown in Lemma 5.5, after repeatedly replacing each (S, T) -stable sub-
 1042 string by a single character, we obtain equivalent strings S' and T' with $|S'| \leq 2k + 1$.
 1043 Under the assumption that $\text{dist}_{TD}(S', T') \leq k$, then each duplication can at most dou-
 1044 ble the size of the previous string. Therefore, T' must have size at most $(2k + 1)2^k$. \square

1045 The kernelization can be performed in polynomial time, as one only needs to
 1046 identify (S, T) -stable pairs and contract them. Running the brute-force algorithm
 1047 from Proposition 5.1 yields the following.

1048 **COROLLARY 5.7.** *The exemplar k -tandem duplication problem can be solved in*
 1049 *time $O(((2k + 1)2^k)^{2k} + \text{poly}(n)) = 2^{O(k^2)} + \text{poly}(n)$, where n is the size of the input.*

1050 **6. Open problems.** Although this work answers some open questions, many of
 1051 them still deserve investigation. We conclude with some of these questions along with
 1052 future research perspectives.

- 1053 • Is the k -TD problem FPT in parameter k ? As we observe in our Exemplar- k -
 1054 TD kernelization, if T and S are large compared to k , they must share many
 1055 long common substrings which could be exploited for an FPT algorithm. It
 1056 is also an interesting question whether Exemplar- k -TD admits a polynomial
 1057 size kernel.

- 1058 • If $|\Sigma| = 2$ or 3 , is the TD problem still NP-hard?
- 1059 • Can one decide in polynomial time whether $S \Rightarrow_* T$? The only known result
- 1060 on this topic is that it can be done if $|\Sigma| = 2$, as one can construct a finite
- 1061 automaton accepting all strings generated by S (though this automaton does
- 1062 not give the minimum number of duplications required).
- 1063 • Does the TD problem admit a constant factor approximation algorithm? The
- 1064 answer might depend on the hardness of deciding whether $S \Rightarrow_* T$, but one
- 1065 might still consider the promise version of the problem.
- 1066 • If the length of each duplicated string is bounded by d , is k -TD in P (with d
- 1067 treated as a constant)? We believe that it is FPT in $k + d$, but is it FPT in
- 1068 d ?
- 1069 • Another parameter of interest is the maximum number of occurrences of a
- 1070 character in S or T . That is, if each $a \in \Sigma$ occurs at most c times in both S
- 1071 and T , is the problem FPT in parameter c ?

1072

REFERENCES

- 1073 [1] N. ALON, J. BRUCK, F. F. HASSANZADEH, AND S. JAIN, *Duplication distance to the root for*
- 1074 *binary sequences*, IEEE Transactions on Information Theory, 63 (2017), pp. 7793–7803,
- 1075 <https://doi.org/10.1109/TIT.2017.2744608>.
- 1076 [2] G. BENSON AND L. DONG, *Reconstructing the duplication history of a tandem repeat*, in Pro-
- 1077 *ceedings of the 17th International Conference on Intelligent Systems for Molecular Biology*,
- 1078 *ISMB'99, AAAI, 1999*, pp. 44–53.
- 1079 [3] D. P. BOVET AND S. VARRICCHIO, *On the regularity of languages on a binary alphabet generated*
- 1080 *by copying systems*, Information Processing Letters, 44 (1992), pp. 119–123, [https://doi.](https://doi.org/10.1016/0020-0190(92)90050-6)
- 1081 [org/10.1016/0020-0190\(92\)90050-6](https://doi.org/10.1016/0020-0190(92)90050-6).
- 1082 [4] L. BULTEAU, G. FERTIN, AND I. RUSU, *Sorting by transpositions is difficult*, SIAM Journal on
- 1083 *Discrete Mathematics*, 26 (2012), pp. 1148–1180, <https://doi.org/10.1137/110851390>.
- 1084 [5] B. CHARLESWORTH, P. SNEGOWSKI, AND W. STEPHAN, *The evolutionary dynamics of repetitive*
- 1085 *dna in eukaryotes*, Nature, 371 (1994), pp. 215–220, <https://doi.org/10.1038/371215a0>.
- 1086 [6] K. CHAUDHURI, K. CHEN, R. MIHAESCU, AND S. RAO, *On the tandem duplication-random*
- 1087 *loss model of genome rearrangement*, in Proceedings of 17th ACM-SIAM Symposium on
- 1088 *Discrete Algorithms, SODA'06, SIAM, 2006*, pp. 564–570.
- 1089 [7] Z.-Z. CHEN, L. WANG, AND Z. WANG, *Approximation algorithms for reconstructing the dupli-*
- 1090 *cation history of tandem repeats*, Algorithmica, 54 (2009), pp. 501–529, [https://doi.org/](https://doi.org/10.1007/s00453-008-9209-8)
- 1091 [10.1007/s00453-008-9209-8](https://doi.org/10.1007/s00453-008-9209-8).
- 1092 [8] D.-J. CHO, Y.-S. HAN, AND H. KIM, *Bound-decreasing duplication system*, Theoretical Com-
- 1093 *puter Science*, 793 (2019), pp. 152–168, <https://doi.org/10.1016/j.tcs.2019.06.018>.
- 1094 [9] J. DASSOW, V. MITRANA, AND G. PAUN, *On the regularity of the duplication closure*, Bulletin
- 1095 *of the EATCS*, 69 (1999), pp. 133–136.
- 1096 [10] R. G. DOWNEY AND M. R. FELLOWS, *Fixed-parameter tractability and completeness ii: On*
- 1097 *completeness for $w[1]$* , Theoretical Computer Science, 141 (1995), pp. 109–131, [https://](https://doi.org/10.1016/0304-3975(94)00097-3)
- 1098 [doi.org/10.1016/0304-3975\(94\)00097-3](https://doi.org/10.1016/0304-3975(94)00097-3).
- 1099 [11] A. EHRENFEUCHT AND G. ROZENBERG, *On regularity of languages generated by copying sys-*
- 1100 *tems*, Discrete Applied Mathematics, 8 (1984), pp. 313–317, [https://doi.org/10.1016/](https://doi.org/10.1016/0166-218X(84)90129-X)
- 1101 [0166-218X\(84\)90129-X](https://doi.org/10.1016/0166-218X(84)90129-X).
- 1102 [12] O. GASCUEL, M. D. HENDY, A. JEAN-MARIE, AND R. MCLACHLAN, *The combinatorics of*
- 1103 *tandem duplication trees*, Systematic Biology, 52 (2003), pp. 110–118, [https://doi.org/10.](https://doi.org/10.1080/10635150390132821)
- 1104 [1080/10635150390132821](https://doi.org/10.1080/10635150390132821).
- 1105 [13] D. GUSFIELD, *Algorithms on strings, trees and sequences: computer science and computational*
- 1106 *biology*, Cambridge University Press, 1997.
- 1107 [14] D. GUSFIELD AND J. STOYE, *Linear time algorithms for finding and representing all the tandem*
- 1108 *repeats in a string*, Journal of Computer and System Sciences, 69 (2004), pp. 525–546,
- 1109 <https://doi.org/10.1016/j.jcss.2004.03.004>.
- 1110 [15] S. HANNENHALI AND P. A. PEVZNER, *Transforming men into mice (polynomial algorithm*
- 1111 *for genomic distance problem)*, in Proceedings of IEEE 36th Symposium on Foundations
- 1112 *of Computer Science, FOCS'95, IEEE, 1995*, pp. 581–592, [https://doi.org/10.1109/SFCS.](https://doi.org/10.1109/SFCS.1995.492588)
- 1113 [1995.492588](https://doi.org/10.1109/SFCS.1995.492588).

- 1114 [16] F. F. HASSANZADEH, M. SCHWARTZ, AND J. BRUCK, *The capacity of string-duplication systems*,
1115 IEEE Transactions on Information Theory, 62 (2016), pp. 811–824, [https://doi.org/10.](https://doi.org/10.1109/TIT.2015.2505735)
1116 [1109/TIT.2015.2505735](https://doi.org/10.1109/TIT.2015.2505735).
- 1117 [17] M. ITO, P. LEUPOLD, AND K. SHIKISHIMA-TSUJI, *Closure of languages under bounded duplica-*
1118 *tions*, in Proceedings of the 10th International Conference on Developments in Language
1119 Theory, vol. 4036 of Lecture Notes in Computer Science, Springer-Verlag, 2006, pp. 238–247,
1120 https://doi.org/10.1007/11779148_22.
- 1121 [18] S. JAIN, F. F. HASSANZADEH, AND J. BRUCK, *Capacity and expressiveness of genomic tandem*
1122 *duplication*, IEEE Transactions on Information Theory, 63 (2017), pp. 6129–6138, [https:](https://doi.org/10.1109/TIT.2017.2728079)
1123 [//doi.org/10.1109/TIT.2017.2728079](https://doi.org/10.1109/TIT.2017.2728079).
- 1124 [19] R. M. KARP, *Reducibility among combinatorial problems*, in Complexity of Computer Com-
1125 putations, R. E. Miller, J. W. Thatcher, and J. D. Bohlinger, eds., The IBM Re-
1126 search Symposia Series, Springer, Boston, MA, 1972, pp. 85–103, [https://doi.org/10.1007/](https://doi.org/10.1007/978-1-4684-2001-2_9)
1127 [978-1-4684-2001-2_9](https://doi.org/10.1007/978-1-4684-2001-2_9).
- 1128 [20] M. LAFOND, B. ZHU, AND P. ZOU, *The tandem duplication distance is NP-hard*, in Proceedings
1129 of STACS’2020, LIPIcs, vol. 154, Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2020,
1130 pp. 15:1–15:15.
- 1131 [21] G. M. LANDAU, J. P. SCHMIDT, AND D. SOKOL, *An algorithm for approximate tandem re-*
1132 *peats*, Journal of Computational Biology, 8 (2001), pp. 1–18, [https://doi.org/10.1089/](https://doi.org/10.1089/106652701300099038)
1133 [106652701300099038](https://doi.org/10.1089/106652701300099038).
- 1134 [22] J. LEECH, *A problem on strings of beads*, The Mathematical Gazette, 41 (1957), pp. 277–278.
- 1135 [23] I. LETUNIC, R. R. COPLEY, AND P. BORK, *Common exon duplication in animals and its*
1136 *role in alternative splicing*, Human Molecular Genetics, 11 (2002), pp. 1561–1567, [https:](https://doi.org/10.1093/hmg/11.13.1561)
1137 [//doi.org/10.1093/hmg/11.13.1561](https://doi.org/10.1093/hmg/11.13.1561).
- 1138 [24] P. LEUPOLD, C. MARTIN-VIDE, AND V. MITRANA, *Uniformly bounded duplication languages*,
1139 Discrete Applied Mathematics, 146 (2005), pp. 301–310, [https://doi.org/10.1016/j.dam.](https://doi.org/10.1016/j.dam.2004.10.003)
1140 [2004.10.003](https://doi.org/10.1016/j.dam.2004.10.003).
- 1141 [25] P. LEUPOLD, V. MITRANA, AND J. M. SEMPERE, *Formal languages arising from gene repeated*
1142 *duplication*, in Aspects of Molecular Computing, G. P. Natasa Jonoska and G. Rozen-
1143 berg, eds., vol. 2950 of Lecture Notes in Computer Science, Springer-Verlag, Berlin, 2004,
1144 pp. 297–308, https://doi.org/10.1007/978-3-540-24635-0_22.
- 1145 [26] M. MACDONALD, ET AL., AND P. S. HARPER, *A novel gene containing a trinucleotide repeat*
1146 *that is expanded and unstable on huntington’s disease*, Cell, 72 (1993), pp. 971–983, [https:](https://doi.org/10.1016/0092-8674(93)90585-E)
1147 [//doi.org/10.1016/0092-8674\(93\)90585-E](https://doi.org/10.1016/0092-8674(93)90585-E).
- 1148 [27] L. OESPER, A. M. RITZ, S. J. AERNI, R. DREBIN, AND B. J. RAPHAEL, *Reconstructing cancer*
1149 *genomes from paired-end sequencing data*, BMC Bioinformatics, 13 (2012), p. S10, [https:](https://doi.org/10.1186/1471-2105-13-S6-S10)
1150 [//doi.org/10.1186/1471-2105-13-S6-S10](https://doi.org/10.1186/1471-2105-13-S6-S10).
- 1151 [28] L. QINGGE, X. HE, Z. LIU, AND B. ZHU, *On the minimum copy number generation problem*
1152 *in cancer genomics*, in Proceedings of the 10th ACM Conference on Bioinformatics, Com-
1153 putational Biology, and Health Informatics, BCB’18, New York, NY, 2018, ACM Press,
1154 pp. 260–269, <https://doi.org/10.1145/3233547.3233586>.
- 1155 [29] D. SANKOFF, *Gene and genome duplication*, Current opinion in genetics and development, 11
1156 (2001), pp. 681–684, [https://doi.org/10.1016/S0959-437X\(00\)00253-7](https://doi.org/10.1016/S0959-437X(00)00253-7).
- 1157 [30] A. J. SHARP, ET AL., AND E. E. EICHLER, *Segmental duplications and copy-number variation*
1158 *in the human genome*, The American J. of Human Genetics, 77 (2005), pp. 78–88, [https:](https://doi.org/10.1086/431652)
1159 [//doi.org/10.1086/431652](https://doi.org/10.1086/431652).
- 1160 [31] J. W. SZOSTAK AND R. WU, *Unequal crossing over in the ribosomal dna of saccharomyces*
1161 *cerevisiae*, Nature, 284 (1980), pp. 426–430, <https://doi.org/10.1038/284426a0>.
- 1162 [32] O. TREMBLAY-SAVARD, D. BERTRAND, AND N. EL-MABROUK, *Evolution of orthologous*
1163 *tandemly arrayed gene clusters*, BMC Bioinformatics, 12 (2011), p. S2, [https://doi.org/10.](https://doi.org/10.1186/1471-2105-12-S9-S2)
1164 [1186/1471-2105-12-S9-S2](https://doi.org/10.1186/1471-2105-12-S9-S2).
- 1165 [33] M.-W. WANG, *On the irregularity of the duplication closure*, Bulletin of the EATCS, 70 (2000),
1166 pp. 162–163.
- 1167 [34] B. ZHU, *Tandem duplication, segmental duplications and deletions, and their applications*, in
1168 Proceedings of the 15th International Computer Science Symposium in Russia (CSR’20),
1169 LNCS, vol. 12159, Springer, 2020, pp. 79–102.