

Revisiting the Parameterized Complexity of Maximum-Duo Preservation String Mapping*

Christian Komusiewicz¹, Mateus de Oliveira Oliveira², and Meirav Zehavi²

1 Friedrich-Schiller-Universität Jena, Germany,
christian.komusiewicz@uni-jena.de

2 University of Bergen, Norway, {mateus.oliveira|meirav.zehavi}@ii.uib.no

Abstract

In the MAXIMUM-DUO PRESERVATION STRING MAPPING (MAX-DUO PSM) problem, the input consists of two related strings A and B of length n and a nonnegative integer k . The objective is to determine whether there exists a mapping m from the set of positions of A to the set of positions of B that maps only to positions with the same character and preserves at least k *duos*, which are pairs of adjacent positions. We develop a randomized algorithm that solves MAX-DUO PSM in time $4^k \cdot n^{O(1)}$, and a deterministic algorithm that solves this problem in time $6.855^k \cdot n^{O(1)}$. The previous best known (deterministic) algorithm for this problem has running time $(8e)^{2k+o(k)} \cdot n^{O(1)}$ [Beretta et al., Theor. Comput. Sci. 2016]. We also show that MAX-DUO PSM admits a problem kernel of size $O(k^3)$, improving upon the previous best known problem kernel of size $O(k^6)$.

1998 ACM Subject Classification G.2.1 Combinatorics, F.2 Analysis of Algorithms and Problem Complexity

Keywords and phrases comparative genomics, parameterized complexity, kernelization

Digital Object Identifier 10.4230/LIPIcs.CPM.2017.44

1 Introduction

Computing distances between strings is a fundamental task in computer science. For many distance measures, the distance between two strings A and B is defined as the minimum number of local operations that are needed to transform A into B , for example the deletion or insertion of a character. For these measures, the distance between two strings A and B can be usually computed in polynomial time [13, 22]. In some applications, however, it is necessary to consider nonlocal operations that transform one string into the other. In comparative genomics, for example, genomes are modeled as strings with one character corresponding to a complete gene and one is interested in determining the evolutionary distance between two genomes. During biological evolution, genomes may be altered by large-scale mutations such as the reversal or the transposition of larger parts of the genome [19].

One approach to approximate the distance between two strings A and B with respect to many of these operations is to compute a smallest *common string partition* [11, 26]. Informally, a size- ℓ common string partition of two strings A and B is a partition of A and B , each into ℓ nonoverlapping substrings, such that the resulting two multisets of substrings

* Christian Komusiewicz gratefully acknowledges support from the DFG, project MAGZ (KO 3669/4-1). Mateus de Oliveira Oliveira gratefully acknowledges support from the Bergen Research Foundation.



© Christian Komusiewicz, Mateus de Oliveira Oliveira and Meirav Zehavi;
licensed under Creative Commons License CC-BY

28th Annual Symposium on Combinatorial Pattern Matching (CPM 2017).

Editors: Juha Kärkkäinen, Jakub Radoszewski, and Wojciech Rytter; Article No. 44; pp. 44:1–44:18

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

of A and B are the same. The problem to compute a smallest common string partition, known as MINIMUM COMMON STRING PARTITION, is NP-hard [11, 21].

An alternative way of defining such a partition is to ask for a partition of A into ℓ nonoverlapping substrings such that permuting the order of these substrings and concatenating them subsequently gives the string B . This second view implies a mapping m that (bijectively) maps each position i of A to a position $m(i)$ of B such that $A[i] = B[m(i)]$. The size of the common string partition is then exactly the number of pairs of consecutive positions i and $i + 1$ such that $m(i) + 1 \neq m(i + 1)$; these positions are called *duos*. Therefore, computing a mapping m that maps only positions with the same characters to each other and maximizes the number k of consecutive positions for which $m(i) + 1 = m(i + 1)$ directly yields a minimum common string partition of A and B . The problem of computing such a mapping is known as MAXIMUM-DUO PRESERVATION STRING MAPPING (MAX-DUO PSM). Since MAX-DUO PSM is simply a dual of the MINIMUM COMMON STRING PARTITION problem, it is NP-hard as well. Motivated by this hardness, we study MAX-DUO PSM from the viewpoint of parameterized algorithmics. More precisely, our aim is to obtain efficient algorithms when the parameter is k , the number of preserved duos. Before describing previous and our results, we give a formal problem definition.

Formal Problem Definition. Let A and B be two strings over a finite set of symbols Σ . Throughout this work, we assume that $|A| = |B| = n$ and that A and B are *related*, that is, B is a permutation of A . A *mapping of A into B* is a (bijective) function $m : [n] \rightarrow [n]$ where for each $i \in [n]$,¹ $A[i] = B[m(i)]$. A *duo* in A is a pair of consecutive positions $(i, i + 1)$ of A . We say that a mapping m *preserves a duo* $(i, i + 1)$ if $m(i) + 1 = m(i + 1)$. Accordingly, the MAX-DUO PSM problem is defined as follows.

MAXIMUM-DUO PRESERVATION STRING MAPPING (MAX-DUO PSM)

Input: Two related strings, A and B , and a nonnegative integer k .

Question: Does there exist a (bijective) mapping m of A into B such that the number of preserved duos is at least k ?

Previous Work. Initially, MAX-DUO PSM has been proposed as an alternative possibility of achieving approximation algorithms for MINIMUM COMMON STRING PARTITION (MCSP) [10], because the best known polynomial-time approximation algorithm has an approximation factor of $O(\log n \log^* n)$ [12]. Consequently, most work on MAX-DUO PSM focuses on approximation algorithms with the first constant-factor approximation algorithm achieving an approximation factor of 4 [6]. This was subsequently improved to a factor of 3.5 [5] and then to a factor of 3.25 [7]. Recently further progress concerning the approximation factor has been reported [18, 27].

Bretta et al. [2, 1] initiated the study of MAX-DUO PSM from the viewpoint of parameterized algorithmics. They studied both the fixed-parameter tractability and the kernelization complexity of MAX-DUO PSM, showing that this problem can be solved in time $(8e)^{2k+o(k)} \cdot n^{O(1)}$, and that it admits a kernel of size $O(k^6)$. Thus, Bretta et al. [2, 1] were the first to show that MAX-DUO PSM is FPT and that it admits a polynomial kernel. The fixed-parameter algorithm of Bretta et al. [2, 1] is based on a combination of color coding and dynamic programming.

¹ We use $[n]$ as shorthand for $\{1, 2, \dots, n\}$.

In comparison with MAX-DUO PSM, MCSP has been investigated more thoroughly from the viewpoint of parameterized algorithms. Damaschke [15] presented the first fixed-parameter algorithms for MCSP, for combined parameters such as “partition size ℓ plus repetition number of the input strings”.² Subsequently, MCSP was shown to be fixed-parameter tractable with the single parameter partition size ℓ [9]. Jiang et al. [23] considered the combined parameter “partition size ℓ plus maximum occurrence d of any character” and showed that MCSP can be solved in time $(d!)^k \cdot n^{O(1)}$. Subsequently, this running time was improved to $O(d^{2k} \cdot kn)$ [8].

Our Contribution. We make two main contributions. First, we develop two algorithms for the MAX-DUO PSM problem that are substantially faster than the (deterministic) algorithm by Bretta et al. [2, 1], which runs in time $(8e)^{2k+o(k)} \cdot n^{O(1)}$. Specifically, we develop a randomized algorithm that solves MAX-DUO PSM in time $4^k \cdot n^{O(1)}$, as well as a deterministic algorithm that solves this problem in time $6.855^k \cdot n^{O(1)}$. Here, in the context of our randomized algorithm, we mean that if we determine that the input is a yes-instance, then this answer is necessarily correct, and if we determine that the input is a no-instance, then this answer is correct with probability at least $9/10$.³ For the purpose of developing our algorithms, we present a reduction from MAX-DUO PSM to a problem of finding paths in an edge-colored graph, which might be of independent interest. This reduction lies at the heart of our algorithms, since by employing advanced tools from the field of parameterized algorithmics, namely, the methods of narrow sieves [4, 3] and representative sets [20], it is possible to quickly solve the resulting graph problem.

Second, we prove that MAX-DUO PSM admits a kernel of size $O(k^3)$, improving upon the kernel of size $O(k^6)$ by Bretta et al. [2].

Preliminaries. We use $[i, j]$ to denote the set $\{i, i + 1, \dots, j\}$ of natural numbers between i and j . Moreover, given a string A , we denote the substring starting at position i and ending at position j by $A[i, j]$. For a (directed) graph G , let $V(G)$ denote the vertex set of G and $E(G)$ the edge set of G .

The field of parameterized algorithmics studies *parameterized problems*, where each problem instance is associated with a *parameter* k , usually a nonnegative integer. Given a parameterized problem, the first question is whether the problem is *fixed-parameter tractable (FPT)*, that is, whether it can be solved in time $f(k) \cdot |X|^{O(1)}$, where f is an arbitrary function that depends *only* on k and $|X|$ is the size of the input instance. In other words, the notion of FPT signifies that the combinatorial explosion can be confined to the parameter k . A second question is whether the problem also admits a *polynomial kernelization*. Here, a problem Π is said to admit a polynomial kernelization if there exists a polynomial-time algorithm that, given an instance (X, k) of Π , outputs an equivalent instance $(\widehat{X}, \widehat{k})$ of Π , called a *kernel*, where $|\widehat{X}| = \widehat{k}^{O(1)}$ and $\widehat{k} \leq k$; kernelization is a mathematical concept that aims to analyze preprocessing procedures in a formal, rigorous manner. For further details, refer to [17, 14].

Due to lack of space, several proofs are deferred to an appendix.

² The repetition number of a nonempty string x is defined as the largest i such that $x = uv^i w$ where v is nonempty.

³ Clearly, the probability of success can be improved by running the algorithm multiple times and determining that the input is a yes-instance if and only if at least one of the calls determined so.

2 Reduction to a Path Finding Problem

In this section, we present a reduction from MAX-DUO PSM to the following graph problem.

LONG BLUE PATH

Input: A directed acyclic graph (DAG) G , an edge-coloring $c : E(G) \rightarrow \{R, B\}$, a vertex-labeling $\ell : V(G) \rightarrow \mathbb{N}$, and nonnegative integers k and r .

Question: Does G contain a directed path P such that

- $|V(P)| \leq r$,
- for all $u, v \in V(P)$, $\ell(u) \neq \ell(v)$, and
- $|\{e \in E(P) : c(e) = B\}| \geq k$.

Construction. Let (A, B, k) be an instance of MAX-DUO PSM. We construct an instance (G, c, ℓ, k, r) of LONG BLUE PATH as follows (here, the parameter k is the same). First, we initialize G to be an empty graph. Now, for every pair of substrings $A[i, j]$ of A and $B[p, q]$ of B such that $j - i \leq k$ and $A[i, j] = B[p, q]$, we insert a directed path $P_{i,j,p,q}$ on $j - i + 1$ new vertices into G whose edges are colored blue and such that the label of the d th vertex on this path is $(p + d - 1)$. The purpose of this path is to represent the possibility to preserve all duos in $A[i, j]$ by mapping this substring to $B[p, q]$. The labels of the vertices are meant to ensure that every position in B is mapped only once. Now, a complete mapping of A to B can be seen as a combination of mappings of substrings that are represented by the paths. Thus, we next turn to connect the paths we have just constructed by adding new edges.

For every two paths $P_{i,j,p,q}$ and $P_{i',j',p',q'}$ such that $j < i'$, we add a red edge from the last vertex of the path $P_{i,j,p,q}$ to the first vertex of the path $P_{i',j',p',q'}$. Informally, the manner in which we direct these edges is meant to ensure that every position in A is mapped only once. Clearly, the resulting graph G is a DAG. Finally, we set $r = 2k$.

Correctness. We first note that the construction can be done in time $O(|V(G)| + |E(G)|)$. Now, observe that the number of paths $P_{i,j,p,q}$ that G contains is bounded by $n^2(k + 1)$ (as the index q equals $p + (j - i)$), and that each path $P_{i,j,p,q}$ consists of at most $(k + 1)$ vertices. Hence, it holds that $|V(G)| \leq n^2(k + 1)^2$ which directly implies $|E(G)| < n^4(k + 1)^2$. Thus, we have the following observation.

► **Observation 1.** *The instance (G, c, ℓ, k, r) can be constructed in time $O(n^4k^2)$.*

We prove the correctness by proving two lemmata that together imply that the instances (A, B, k) and (G, c, ℓ, k, r) are equivalent.

► **Lemma 1.** *If (A, B, k) is a yes-instance of MAX-DUO PSM, then (G, c, ℓ, k, r) is a yes-instance of LONG BLUE PATH.*

Proof. Let m be a mapping from A into B preserving at least k duos. Consider the set $\{A_1, \dots, A_r\}$ of substrings of A containing exactly the first k preserved duos, where we assume that A_i precedes A_{i+1} in A . Consider any A_z and let $[i_z, j_z]$ be the set of positions of A_z in A . Since the mapping preserves the duos in A_z , there is a substring $B[q_z, p_z]$ such that $m(i_z + s) = q_z + s$, $0 \leq s \leq j - i$. This implies that $A[i_z, j_z] = B[q_z, p_z]$. Thus, G contains the path $P_z := P_{i_z, j_z, q_z, p_z}$.

By the above, for each A_z , G contains a path P_z containing $|A_z| - 1$ blue edges. Moreover, for A_i and A_j , the vertices in P_i and P_j have different labels since the mapping m is injective. Finally, there is a red edge from the last vertex of P_i to the first vertex of P_{i+1} since the last position of A_i is strictly smaller than the first position of A_{i+1} . Thus, the concatenation

of P_1, P_2 until P_r gives a path in G . The number of blue edges in this path is exactly k , and the number of vertices in this path is at most $2k$, since every P_i contains at least one blue edge. ◀

► **Lemma 2.** *If (G, c, ℓ, k) is a yes-instance of LONG BLUE PATH, then (A, B, k) is a yes-instance of MAX-DUO PSM.*

Proof. Let P be a solution of the LONG BLUE PATH instance. That is, P is a path in G on at most $2k$ vertices, all with different labels, containing at least k blue edges. Let $\{P_1, \dots, P_r\}$ be the set of disjoint paths obtained from P by removing all red edges where we assume that there is a red edge from P_i to P_{i+1} for all $i \in [r-1]$. Consider some P_z . By the construction of G , $P_z = P_{i,j,q,p}$ for some $i < j$ and $q < p$. Hence, there is a substring $A[i, j]$ of A and a substring $B[q, p]$ of B such that $A[i, j] = B[q, p]$. Call these two substrings the substrings of A and B , respectively, that *correspond* to P_z . Observe that for P_i and P_j , $i < j$, the substrings corresponding to P_i and P_j are disjoint: For the substrings in A this is due to the fact that the indices of the corresponding substring for P_i are lower than those of the substring of A corresponding to P_j . For the substrings in B this is due to the fact that the vertices in P_i and P_j have different labels. Thus, there is a mapping from A into B that maps the corresponding strings for each path P_i and maps all other positions arbitrarily. The number of duos preserved by this mapping is at least k . ◀

Altogether, we arrive at the following.

► **Lemma 3.** *Given an instance (A, B, k) of MAX-DUO PSM, an equivalent instance (G, c, ℓ, k, r) of LONG BLUE PATH where $r = 2k$ can be constructed in time $O(n^4 k^2)$.*

3 A Randomized Algorithm based on Narrow Sieves

In this section, we adapt the method of *narrow sieves* that was applied to solve the k -PATH problem [4] to solve LONG BLUE PATH. More precisely, our objective is to provide a constructive proof for the following result.

► **Lemma 4.** *There exists a randomized algorithm that solves LONG BLUE PATH in time $2^r \cdot r^{O(1)} \cdot |E(G)|$ and polynomial space.*

In light of Lemma 3, once we have Lemma 4 at hand, we immediately obtain the following theorem.

► **Theorem 5.** *There exists a randomized algorithm that solves MAX-DUO PSM in time $4^k \cdot k^{O(1)} \cdot n^4$ and polynomial space.*

In the following, we focus on the proof of Lemma 4. To this end, let (G, c, ℓ, k, r) be an instance of LONG BLUE PATH. Clearly, we can assume that $|V(G)| \leq |E(G)|$. To be able to rely on dynamic programming later, we need to define a notion of a partial solution:

► **Definition 6.** Let P be a directed path in G . Given a vertex $v \in V(G)$, $s \in [r]$ and $b \in [r] \cup \{0\}$, we say that P is a (v, s, b) -path if the last vertex of P is v , $|V(P)| = s$ and $|\{e \in E(P) : c(e) = B\}| = b$. If for all $u, w \in V(P)$, it holds that $\ell(u) \neq \ell(w)$, then we say that P is a *good path*.

To employ the method of narrow sieves, we need to associate labels with entities whose uniqueness should be preserved. For this purpose, we have the following definition:

► **Definition 7.** Let P be a (v, s, b) -path. Given $f : V(P) \rightarrow [r]$, we say that (P, f) is a (v, s, b) -pair. If P is good, then we say that (P, f) is a *good pair*, and if f is an injective function, then we say that (P, f) is an *injective pair*. Given $L \subseteq [r]$ such that the image of f is a subset of L , we say that (P, f) is an L -labeled pair.

Now, we define two central sets of labeled partial solutions. The first one, \mathcal{P} , consists of every pair (P, f) that is an injective (v, s, b) -pair for some $v \in V(G)$ and $s, b \in [r]$ such that $b \geq k$. The second one, \mathcal{Q} , consists of every good pair (P, f) in \mathcal{P} . Note that for every pair $(P, f) \in \mathcal{Q}$, it holds that P is a solution for LONG BLUE PATH, and for every solution P for LONG BLUE PATH, by letting f be a function that assigns i to the i th vertex on P , we obtain a pair $(P, f) \in \mathcal{Q}$. Thus, we have the following observation.

► **Observation 2.** *The instance (G, c, ℓ, k, r) is a yes-instance if and only if $\mathcal{Q} \neq \emptyset$.*

With these definitions at hand, we may describe the rough idea of the approach. We represent all labeled partial solutions of \mathcal{P} by a polynomial in such a way that each labeled partial solution corresponds to one monomial. We will ensure that the partial solutions of $\mathcal{P} \setminus \mathcal{Q}$ cancel each other out which will imply that the polynomial is not identically 0 if and only if $\mathcal{Q} \neq \emptyset$. To this end, we now describe how we represent labeled partial solutions by monomials. For every label $i \in \text{image}(\ell)$ and integer $j \in [r]$, we introduce the variable $x_{i,j}$, and for every edge $e \in E(G)$, we introduce the variable y_e . This gives the following representation:

► **Definition 8.** Let (P, f) be a (v, s, b) -pair. Then, the monomial associated with (P, f) is defined as follows.

$$\text{mon}(P, f) = \prod_{v \in V(P)} x_{\ell(v), f(v)} \cdot \prod_{e \in E(P)} y_e.$$

Accordingly, we define the following polynomial (which would be evaluated over a field of characteristic 2).

► **Definition 9.** $\text{POL} = \sum_{(P, f) \in \mathcal{P}} \text{mon}(P, f).$

To analyze this polynomial, we first observe that given a monomial associated with a pair $(P, f) \in \mathcal{Q}$, we can uniquely recover the pair (P, f) . To see this, consider some monomial M that is associated with a pair $(P, f) \in \mathcal{Q}$. Then, the variables y_e of M specify exactly which edges are used by P , and therefore the path P is recovered. Now, since the pair (P, f) belongs to \mathcal{Q} , we have that P is a good path. Hence, the variables $x_{i,j}$ of M specify exactly how f labels the vertices of P . In other words, we have the following observation.

► **Observation 3.** *For all $(P, f) \in \mathcal{Q}$, there does not exist $(P', f') \in \mathcal{P} \setminus \{(P, f)\}$ such that $\text{mon}(P, f) = \text{mon}(P', f')$.*

The following lemma will be used to show that the partial solutions of $\mathcal{P} \setminus \mathcal{Q}$ cancel each other out.

► **Lemma 10.** *There exists a function $g : \mathcal{P} \setminus \mathcal{Q} \rightarrow \mathcal{P} \setminus \mathcal{Q}$ such that for all $(P, f) \in \mathcal{P} \setminus \mathcal{Q}$, it holds that $\text{mon}(P, f) = \text{mon}(g(P, f))$, $g(P, f) \neq (P, f)$, and $g(g(P, f)) = (P, f)$.*

Proof. Let $<$ be some order on $\{\{u, v\} : u, v \in V(P)\}$. Given $(P, f) \in \mathcal{P} \setminus \mathcal{Q}$, define $\text{rep}(P, f) = \{\{u, v\} : u, v \in V(P), u \neq v, \ell(u) = \ell(v)\}$. Since P is not a good path, it holds that $\text{rep}(P, f) \neq \emptyset$. Hence, it is well defined to let $\{u, v\}$ be the smallest set in $\text{rep}(P, f)$ according to $<$. We let h be defined as f except that $h(u) = f(v)$ and $h(v) = f(u)$. Now, we

set $g(P, f) = (P, h_{P,f})$. Clearly, $g(P, f) \in \mathcal{P}$. Note that $\text{rep}(P, f) = \text{rep}(P, h_{P,f})$, and hence $g(P, f) \notin \mathcal{Q}$ and $g(g(P, f)) = (P, f)$. Since $(P, f) \in \mathcal{P}$, it holds that f is an injective function; therefore $f(v) \neq f(u)$, which implies that $g(P, f) \neq (P, f)$. Finally, since $\ell(u) = \ell(v)$, it holds that $\text{mon}(P, f) = \text{mon}(g(P, f))$. \blacktriangleleft

Let \mathbb{F} be a field of characteristic 2 (to be determined). From now on, we suppose that POL is evaluated over \mathbb{F} . Notice that

$$\text{POL} = \sum_{(P,f) \in \mathcal{Q}} \text{mon}(P, f) + \sum_{(P,f) \in \mathcal{P} \setminus \mathcal{Q}} \text{mon}(P, f).$$

Suppose that POL is evaluated over \mathbb{F} . By Lemma 10, we have that $\text{POL} = \sum_{(P,f) \in \mathcal{Q}} \text{mon}(P, f)$. Then, by Observation 3, we have that POL is not identically 0 if and only if \mathcal{Q} is not empty. Hence, by Observation 2, we have the following lemma.

► **Lemma 11.** *The instance (G, c, ℓ, k, r) is a yes-instance if and only if POL is not identically 0.*

In light of Lemma 11, our task is to determine whether POL is identically 0. For this purpose, we need the following notation. Given $v \in V(G)$, $s \in [r]$, $b \in [r] \cup \{0\}$ and $L \subseteq [r]$, let $\mathcal{P}_{v,s,b,L}$ denote the set of every L -labeled (v, s, b) -pair (P, f) , and $\text{POL}_{v,s,b,L} = \sum_{(P,f) \in \mathcal{P}_{v,s,b,L}} \text{mon}(P, f)$. Moreover, denote

$$\mathcal{P}_L = \bigcup_{v \in V(G), s, b \in [r], b \geq k} \mathcal{P}_{v,s,b,L},$$

and $\text{POL}_L = \sum_{(P,f) \in \mathcal{P}_L} \text{mon}(P, f)$. By the principle of inclusion-exclusion, we have that $\text{POL} = \sum_{L \subseteq [r]} (-1)^{r-|L|} \text{POL}_L$. Then, since \mathbb{F} is a field of characteristic 2 (refer to [4] for further details) we obtain the following.

► **Observation 4.** $\text{POL} = \sum_{L \subseteq [r]} \text{POL}_L$.

Hence, to determine whether POL is identically 0, it is sufficient to determine whether $\sum_{L \subseteq [r]} \text{POL}_L$ is identically 0. To proceed, we need to recall the following well-known lemma.

► **Lemma 12** ([24, 28, 16]). *Let $p(x_1, x_2, \dots, x_n)$ be a nonzero polynomial of total degree at most d over a finite field \mathbb{K} . Then, for $a_1, a_2, \dots, a_n \in \mathbb{K}$ selected independently and uniformly at random, $\Pr(p(a_1, a_2, \dots, a_n) \neq 0) \geq 1 - d/|\mathbb{K}|$.*

Notice that POL is a polynomial of total degree at most $2r$. Therefore, by setting $|\mathbb{F}| = 2^{\lceil \log(20r) \rceil}$, from Lemma 11, Observation 4, and Lemma 12, we have that

► **Lemma 13.** *For a random assignment to all variables $x_{i,j}$ and y_e , if (G, c, ℓ, k, r) is a no-instance, then $\sum_{L \subseteq [r]} \text{POL}_L$ evaluates to 0, and otherwise it does not evaluate to a 0 with probability at least $9/10$.*

In light of Lemma 13, to conclude that Lemma 4 is correct, it is sufficient to prove the following result.

► **Lemma 14.** *Given $L \subseteq [r]$ and an assignment to all variables $x_{i,j}$ and y_e , the polynomial POL_L can be evaluated in time $r^{O(1)} \cdot |E(G)|$.*

Finally, we would like to remark that if one is interested in finding a mapping that is a solution rather than just determining whether such a mapping exists, this goal can be achieved by standard means of self-reduction. Briefly, if k is not positive, then we are done. Else, if the algorithm determines that there exists a solution, then we may “guess” (i.e., perform exhaustive search) a longest substring A' of A that is mapped by some solution while preserving all duos in A' as well as the substring B' of B to which it is mapped. If our guess is correct, then the symbol preceding A' in A is not equal to the symbol preceding B' in B and the symbol after A' in A is also not equal to the symbol after B' in B (if such symbols exist). Then, we may replace A' and B' in A and B , respectively, by some new symbol, decrease k by $|A'| - 1$, and call the algorithm recursively. Notice that the length of A' should be at least 2, and hence the size of the input has decreased.

4 Deterministic Algorithm: Representative Sets

In this section, we adapt the approach in which the method of *representative sets* is applied to solve the k -PATH problem [20]. More precisely, our objective is to provide a constructive proof for the following result.

► **Lemma 15.** *There exists a deterministic algorithm that solves LONG BLUE PATH in time $O\left(\left(\frac{1+\sqrt{5}}{2}\right)^{r+o(r)} \cdot |E(G)| \cdot \log |E(G)|\right)$.*

Combining Lemma 3 and 15 gives us the following.

► **Theorem 16.** *There exists a deterministic algorithm that solves MAX-DUO PSM in time $O\left(\left(\frac{1+\sqrt{5}}{2}\right)^{2k+o(k)} \cdot n^4 \log n\right) = O(6.855^k \cdot n^4 \log n)$.*

5 A Cubic Problem Kernel

In this section we will show that MAX-DUO PSM admits a kernel of size $O(k^3)$. Let (A, B, k) be an instance of MAX-DUO PSM, and let $S \in \{A, B\}$. If $S = A$, then we let $\bar{S} = B$. Analogously, if $S = B$, then we let $\bar{S} = A$.

Let m be a map of S into \bar{S} , and let D be a set of duos. We denote by $m(D) = \{(m(i), m(i+1)) \mid (i, i+1) \in D\}$ the image of D under m . We say that m *preserves* D if m preserves each duo in D . Let C_A and C_B be sets of duos. We say that the pair (C_A, C_B) is *complete* for (A, B, k) if whenever there is a map m of A into B that preserves k duos, then there is a subset $D \subseteq C_A$ with $|D| = k$ and a map m' such that m' preserves D and $m'(D) \subseteq C_B$. The size of (C_A, C_B) is defined as $|C_A| + |C_B|$. Let $f : \mathbb{N} \rightarrow \mathbb{N}$ be a function. A complete pair (C_A, C_B) of size $f(k)$ for (A, B, k) can be used to construct a kernel (A', B', k) of size $O(f(k))$ for (A, B, k) .

► **Theorem 17** ([2, Section 4.2]). *Let (C_A, C_B) be a complete pair of size $f(k)$ for (A, B, k) . Then one can construct in time $O(f(k))$ related strings A' and B' , each of size $O(f(k))$ such that (A, B, k) is a yes-instance of Max-Duo PSM if and only if (A', B', k) is a yes-instance of Max-Duo PSM.*

Using Theorem 17, it is sufficient to show that one can obtain in polynomial time a complete pair (C_A, C_B) for (A, B, k) of size $O(k^3)$.

A *block* of size s is a set $X = \{(i, i + 1), (i + 1, i + 2), \dots, (i + s - 1, i + s)\}$ consisting of s consecutive duos. We say that $(i, i + 1)$ is the *root* of X . If S is a string of length at least $i + s$, then we let $\text{str}(S, X) = S[i, i + s]$ be the substring of S corresponding to the positions that occur in X . The following observation is immediate.

► **Observation 5.** *Let (A, B, k) be an instance of MAX-DUO PSM and let m be a map of A into B that preserves a block X of size k . Then (A, B, k) is a yes-instance of MAX-DUO PSM. Additionally, the instance (A', B', k) where $A' = \text{str}(A, X)$ and $B' = \text{str}(B, m(X))$ is also yes-instance of MAX-DUO PSM.*

In the remainder of this section we assume that no map m of A into B preserves a block of size k . Our algorithm is based on the notion of *rare* duo, which we define next. For each two symbols $a, b \in \Sigma$, and each string $S \in \{A, B\}$, we let

$$n(S, a, b) := |\{i : 1 \leq i \leq |S| - 1, S[i, i + 1] = ab\}|$$

be the number of occurrences of the length-two string ab as a substring of S . We say that a length-two string ab is *rare for S* if ab occurs as a sub-string of both S and \bar{S} and $n(S, a, b) \leq n(\bar{S}, a, b)$. Observe that if ab occurs as many times in S as it occurs in \bar{S} , then ab is rare for both S and \bar{S} . We say that a duo $(i, i + 1)$ is rare for S if $S[i, i + 1]$ is rare for S . We let $\text{rare}(S)$ be the set of duos that are rare for S .

► **Lemma 18.** *If either $|\text{rare}(A)| \geq 4k$ or $|\text{rare}(B)| \geq 4k$, then (A, B, k) is a yes-instance.*

Proof. The duo graph associated with A and B is the bipartite graph $G(A, B) = (V_A \cup V_B, E)$ defined as follows.

$$V_A = \{(i, i + 1) \mid 1 \leq i \leq n - 1\} \quad V_B = \{(j, j + 1) \mid 1 \leq j \leq n - 1\}$$

$$E = \{[(i, i + 1), (j, j + 1)] \mid A[i] = B[j], A[i + 1] = B[j + 1]\}$$

Intuitively, each of the sets V_A and V_B contains all pairs of consecutive positions from $[n]$. A duo $(i, i + 1)$ in V_A is connected to a duo $(j, j + 1)$ in V_B if and only if the length-two string $A[i]A[i + 1]$ is equal to $B[j]B[j + 1]$.

If $e = [(i, i + 1), (j, j + 1)]$ is an edge of $G(A, B)$, then we say that $(i, i + 1)$ is the left endpoint of e and $(j, j + 1)$ is the right endpoint of e . If M is a matching in $G(A, B)$, then we let M_A be the set of duos in V_A that are left endpoints of edges in M , and M_B be the set of duos in V_B that are right endpoints of edges in M .

Assume that either $|\text{rare}(A)| \geq 4k$ or $|\text{rare}(B)| \geq 4k$. Then G contains a matching of size at least $4k$. Let \mathbf{M} be a maximum matching in $G(A, B)$.

It has been shown in [6] that given a matching \mathbf{M} of size at least $4k$ for the graph $G(A, B)$, one can construct a sub-matching M of \mathbf{M} of size at least k such that M directly gives a map preserving at least k duos. Therefore, the instance is a yes-instance in this case. ◀

In the remainder of this section we thus assume that there are less than $4k$ duos that are rare for A , and less than $4k$ duos that are rare for B . This implies that we may add all rare duos to the sets C_A and C_B without surpassing the desired size bound of $O(k^3)$.

Let S be a string in $\{A, B\}$. We say that a duo $(j, j + 1)$ is a *match for a duo* $(i, i + 1)$ in \bar{S} if there exists a map m of S into \bar{S} that preserves $(i, i + 1)$, and $(m(i), m(i + 1)) = (j, j + 1)$. If X and Y are blocks, then we say that Y is a *match for X* in \bar{S} if there exists a map m of S into \bar{S} such that m preserves X , and $m(X) = Y$.

Algorithm 1

```

1: procedure ROOTS( $S, i, i + 1$ )
2:    $R = \emptyset$ 
3:    $k' \leftarrow$  size of the maximal block which is rooted at  $(i, i + 1)$ , rare for  $S$ , and has a
4:     match in  $\bar{S}$ . Note that  $k' \leq k - 1$ .
5:   for  $\ell = k'$  to 1 do
6:      $X \leftarrow$  unique block of size  $\ell$  rooted at  $(i, i + 1)$ 
7:     for  $j = 1$  to  $n - 1$  do
8:       if  $|R| < 2k - 1$  and  $|j' - j| > k \forall j' \in R$  and
9:          $(j, j + 1)$  is a root for a match of  $X$  in  $\bar{S}$  then
10:         $R \leftarrow R \cup \{(j, j + 1)\}$ 
11:   output  $R$ 

```

► **Observation 6.** Let $S \in \{A, B\}$ and let $(j, j + 1)$ be a match for $(i, i + 1)$ in \bar{S} . Then if $(i, i + 1)$ is not rare for S , $(j, j + 1)$ is rare for \bar{S} .

Proof. Since $(j, j + 1)$ is a match for $(i, i + 1)$ in \bar{S} , there is some length-two string ab such that $S[i]S[i + 1] = \bar{S}[j]\bar{S}[j + 1] = ab$. Since $(i, i + 1)$ is not rare for S , the string ab occurs strictly more often in S than it occurs in \bar{S} . In other words, $n(S, a, b) > n(\bar{S}, a, b)$. This implies that $(j, j + 1)$ is rare for \bar{S} . ◀

This observation is useful because it tells us that for each match in a map, one of the two duos is rare, so by adding all the rare duos to C_A and C_B , we essentially pick up one half of each match. We now consider two types of matched blocks that may occur in the solution. First, there may be pairs of matched blocks X and Y that both contain nonrare duos. We can add all duos of these blocks by considering a sufficiently large neighborhood of all rare duos. To this end, for each $i \in \{1, \dots, n - 1\}$, let

$$\mathcal{B}_k(i) = \{(i', i' + 1) \mid i' \in \{1, \dots, n - 1\}, i - k \leq i' \leq i + k\}$$

denote the *ball of radius k* around the duo $(i, i + 1)$

The following lemma essentially implies that by adding the ball of radius k around each rare duo, we add all pairs of matched blocks that both contain at least one nonrare duo.

► **Lemma 19.** Let $S \in \{A, B\}$, X be a block of size at most $k - 1$ containing a duo $(i, i + 1)$ that is not rare for S , and let m be a map of S into \bar{S} such that X is preserved by m . Then $(m(i), m(i + 1))$ is rare for \bar{S} and $m(X) \subseteq \mathcal{B}_k(m(i))$.

Proof. Since $(i, i + 1)$ is preserved by m , $(m(i), m(i + 1))$ is a match for $(i, i + 1)$ in \bar{S} . Since $(i, i + 1)$ is not rare for S , by Observation 6, $(m(i), m(i + 1))$ is rare for \bar{S} . Since m preserves X and since $|X| \leq k - 1$, $m(X)$ is a block of size at most $k - 1$. Therefore, all duos in $m(X)$ must be in the ball of radius k around $(m(i), m(i + 1))$, that is, $m(X) \subseteq \mathcal{B}_k(m(i))$. ◀

We now turn to the second type of matched pairs of blocks, those where one block X of S has only rare duos for S ; we call such a block X *rare*. Since X is rare, it is rooted at some rare duo $(i, i + 1)$. To obtain the complete set, we need to add duos in \bar{S} . This is done by the procedure ROOTS which receives as input a string S and a duo in S and returns a set of duos ROOTS($S, i, i + 1$).

Intuitively, for each block X that is rare for S and rooted at $(i, i + 1)$, the set ROOTS($S, i, i + 1$) contains a selection of roots of matches for X in the string \bar{S} . This selection is made

according to two criteria. First, roots of matches for larger blocks are added first. Second, the roots in $\text{ROOTS}(S, i, i+1)$ are sufficiently far apart from each other. Now consider the set

$$F(S, i, i+1) = \bigcup_{(j, j+1) \in \text{ROOTS}(S, i, i+1)} \mathcal{B}_k(j).$$

Intuitively, $F(S, i, i+1)$ consists of all duos that are sufficiently close to duos in $\text{ROOTS}(S, i, i+1)$. The next lemma states that if some map m of S into \bar{S} preserves some block X that is rooted at $(i, i+1)$ and rare for S , then this map can be transformed into a map m' that preserves X , that sends X to $F(S, i, i+1)$, and that is equal to m on every duo outside X .

► **Lemma 20.** *Let m be a map of S into \bar{S} , D be a set of duos such that $|D| = k$, and $X \subseteq D$ be a block that is rooted at $(i, i+1)$, that is rare for S and that is preserved by m . Then there is a map m' of S into \bar{S} such that X is preserved by m' , such that $m'(X) \subseteq F(S, i, i+1)$ and such that $(m'(i'), m'(i'+1)) = (m(i'), m(i'+1))$ for each $(i', i'+1) \in D \setminus X$.*

Proof. Let $(j, j+1)$ be the root of $m(X)$ in \bar{S} . Let $n(S, X)$ be the number of duos in $\text{ROOTS}(S, i, i+1)$ that are roots of matches for X in \bar{S} . Suppose that $n(S, X) < 2k - 1$. Then either $(j, j+1) \in \text{ROOTS}(S, i, i+1)$ or $(j, j+1)$ does not belong to $\text{ROOTS}(S, i, i+1)$ and there exists some duo $(j', j'+1) \in \text{ROOTS}(S, i, i+1)$ with $|j' - j| < k$. Note that if this were not the case, the duo $(j, j+1)$ would have been added to $\text{ROOTS}(S, i, i+1)$, since all three conditions of the 'If' instruction of Algorithm 1 would have been satisfied. In any case, $m(X) \subseteq \mathcal{B}_k(j') \subseteq F(S, i, i+1)$. Therefore, if $n(S, X) < 2k - 1$, we may simply set $m' = m$.

Now assume that $n(S, X) = 2k - 1$. Note that for each duo $(j, j+1)$ there are no three distinct j_1, j_2 and j_3 such that $(j_l, j_l+1) \in \text{ROOTS}(S, i, i+1)$ and $(j, j+1) \in \mathcal{B}(j_l)$ for $l \in \{1, 2, 3\}$. In other words $(j, j+1)$ can intersect at most two balls of radius k rooted at duos in $\text{ROOTS}(S, i, i+1)$. Therefore, since $|D \setminus X| \leq k - 1$, the set $D \setminus X$ intersects at most $2k - 2$ balls of radius r rooted at duos in $\text{ROOTS}(S, i, i+1)$. In other words, there is at least one $(j', j'+1) \in \text{ROOTS}(S, i, i+1)$ that is the root of a match for X in S and such that $\mathcal{B}_k(j') \cap (D \setminus X) = \emptyset$. Therefore, we may set m' as the map of S into \bar{S} that preserves X , that sends the root of X to $(j', j'+1)$, and that is equal to m on every duo $(i', i'+1) \in D \setminus X$. ◀

Now, for each $S \in \{A, B\}$, consider the following set C_S of duos.

$$C_S = \left[\bigcup_{(i, i+1) \in \text{rare}(S)} \mathcal{B}_k(i) \right] \cup \left[\bigcup_{(i, i+1) \in \text{rare}(\bar{S})} F(\bar{S}, i, i+1) \right]. \quad (1)$$

In other words, for each duo $(i, i+1)$ that is rare for S , C_S contains all duos in the ball of radius k around $(i, i+1)$. Moreover, for each duo $(i, i+1)$ that is rare for \bar{S} , C_S contains all duos in the set $F(\bar{S}, i, i+1)$. The following lemma states that if a map m of S into \bar{S} preserves a set D containing k duos, then there exists a map m' that also preserves D in such a way that $m'(D) \subseteq C_{\bar{S}}$.

► **Lemma 21.** *Let D be a set of duos such that $|D| = k$. Let m be a map of S into \bar{S} that preserves all duos in D . Then there is a map m' of S into \bar{S} that preserves all duos in D , and such that $m'(D) \subseteq C_{\bar{S}}$.*

Proof. Let X_1, \dots, X_r be the set of rare blocks that are contained in D and that are maximal with respect to set inclusion. In other words, for each $j \in \{1, \dots, r\}$ and each Y such that $X_j \subseteq Y \subseteq D$, we have that Y is not a rare block. Note that since these blocks are rare and

maximal, they are pairwise disjoint, i.e., $X_j \cap X_{j'} = \emptyset$ for $j \neq j'$. For each $j \in \{1, \dots, r\}$ let $(i_j, i_j + 1)$ be the root of X_j and $D_j = D \setminus X_j$. Additionally, let $D' = D \setminus \bigcup_{j=1}^r X_j$. Note that $D' \subseteq D_j$ for each $j \in \{1, \dots, r\}$.

Let m_0, m_1, \dots, m_r be maps of S into \bar{S} defined inductively as follows. First, we set $m_0 = m$. Now, for each $j \in \{1, \dots, r\}$, we let m_j be a map of S into \bar{S} constructed according to Lemma 20. More precisely, m_j preserves X_j , $m_j(X_j) \subseteq F(S, i_j, i_j + 1)$, and $(m_j(i), m_j(i + 1)) = (m_{j-1}(i), m_{j-1}(i + 1))$ for each duo $(i, i + 1) \in D_j = D \setminus X_j$.

Using the maps m_0, m_1, \dots, m_r defined above, it follows by induction on j that for each $l \in \{1, \dots, j\}$, $m_j(X_l) \subseteq F(S, i_l, i_l + 1) \subseteq C_{\bar{S}}$ and $(m_j(i), m_j(i + 1)) = (m(i), m(i + 1))$ for each $(i, i + 1) \in D'$. In particular, for each $l \in \{1, \dots, r\}$, $m_r(X_l) \subseteq F(S, i_l, i_l + 1) \subseteq C_{\bar{S}}$ and $(m_r(i), m_r(i + 1)) = (m(i), m(i + 1))$ for each $(i, i + 1) \in D' \subseteq D_j$. This shows, that m_r preserves D , and sends $\bigcup_{j=1}^r X_j$ to a subset of $C_{\bar{S}}$ and agrees with m in every duo in $D' = D \setminus \bigcup_{j=1}^r X_j$.

Let $m' = m_r$. It remains to show that m' also sends blocks that are not rare for S to subsets of $C_{\bar{S}}$. Let X'_1, \dots, X'_s be the maximal blocks that are contained in D and that are not rare for S . Note that these blocks are indeed contained in D' and form a partition of D' . Since for each $j \in \{1, \dots, s\}$, X'_j has at least one duo $(i, i + 1)$ that is not rare for S , Lemma 19 implies that $(m'(i), m'(i + 1))$ is rare for \bar{S} and that $m'(X'_j) \subseteq \mathcal{B}_k(m'(i)) \subseteq C_{\bar{S}}$. Since X'_1, \dots, X'_s forms a partition of D' , $m'(D') \subseteq C_{\bar{S}}$. Since by the discussion above, $m'(\bigcup_{j=1}^r X_j)$ is also a subset of $C_{\bar{S}}$, we have that $m'(D) \subseteq C_{\bar{S}}$. ◀

Let C_A and C_B be sets of duos constructed according to Equation 1. We can show that (C_A, C_B) is complete for (A, B, k) by applying Lemma 21 twice. More precisely, once with respect to maps of A into B , and once with respect to maps of B into A .

► **Lemma 22.** *The pair (C_A, C_B) is complete for (A, B, k) .*

Proof. Let D_1 be a set of duos of size k . Let m_1 be a map of A into B which preserves all duos in D_1 . Then by Lemma 21 there is a map m_2 of A into B which also preserves all duos in D_1 , but with the property that $m_2(D_1) \subseteq C_B$. Now let $D_2 = m_2(D_1)$, and $m_3 = m_2^{-1}$ be the inverse of m_2 . In other words, m_3 is a map of B into A such that for each $i \in [n]$, $m_2(i) = j$ if and only if $m_3(j) = i$. Then m_3 preserves all duos in D_2 . By Lemma 21 there is a map m_4 of B into A that also preserves all duos in D_2 but with the additional property that $m_4(D_2) \subseteq C_A$.

Let $D_3 = m_4(D_2)$, and let $m_5 = m_4^{-1}$ be the inverse of m_4 . Then m_5 is a map of A into B that preserves $D_3 \subseteq C_A$ and such that $m_5(D_3) = D_2 \subseteq C_B$. Since $|D_3| = |D_2| = k$, the pair (C_A, C_B) is complete for (A, B, k) . ◀

Now, we can upper-bound the size of C_S and the time needed to construct C_S , thus arriving at our main theorem.

► **Theorem 23.** *Given an instance $I = (A, B, k)$ of MAX-DUO PSM, one can construct in time $O(|\Sigma|^2 \cdot n + k^3 \cdot n)$ an instance $I' = (A', B', k)$ of MAX-DUO PSM with $|A'|$ and $|B'|$ bounded by $O(k^3)$ such that I is a yes-instance if and only if I' is a yes-instance.*

References

- 1 Stefano Beretta, Mauro Castelli, and Riccardo Dondi. Corrigendum to “Parameterized tractability of the maximum-duo preservation string mapping problem” [Theoret. Comput. Sci. 646 (2016) 16–25]. *Theor. Comput. Sci.*, 653:108–10, 2016.

- 2 Stefano Beretta, Mauro Castelli, and Riccardo Dondi. Parameterized tractability of the maximum-duo preservation string mapping problem. *Theor. Comput. Sci.*, 646:16–25, 2016.
- 3 Andreas Björklund. Determinant sums for undirected hamiltonicity. *SIAM J. Comput.*, 43(1):280–299, 2014.
- 4 Andreas Björklund, Thore Husfeldt, Petteri Kaski, and Mikko Koivisto. Narrow sieves for parameterized paths and packings. *CoRR*, abs/1007.1161, 2010.
- 5 Nicolas Boria, Gianpiero Cabodi, Paolo Camurati, Marco Palena, Paolo Pasini, and Stefano Quer. A $7/2$ -approximation algorithm for the maximum duo-preservation string mapping problem. In Roberto Grossi and Moshe Lewenstein, editors, *Proceedings of the 27th Annual Symposium on Combinatorial Pattern Matching (CPM '16)*, volume 54 of *LIPICs*, pages 11:1–11:8. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 2016.
- 6 Nicolas Boria, Adam Kurpisz, Samuli Leppänen, and Monaldo Mastrolilli. Improved approximation for the maximum duo-preservation string mapping problem. In Daniel G. Brown and Burkhard Morgenstern, editors, *Proceedings of the 14th International Workshop on Algorithms in Bioinformatics (WABI '14)*, volume 8701 of *Lecture Notes in Computer Science*, pages 14–25. Springer, 2014.
- 7 Brian Brubach. Further improvement in approximating the maximum duo-preservation string mapping problem. In Martin C. Frith and Christian Nørgaard Storm Pedersen, editors, *Proceedings of the 16th International Workshop on Algorithms in Bioinformatics (WABI '16)*, volume 9838 of *Lecture Notes in Computer Science*, pages 52–64. Springer, 2016.
- 8 Laurent Bulteau, Guillaume Fertin, Christian Komusiewicz, and Irena Rusu. A fixed-parameter algorithm for minimum common string partition with few duplications. In Aaron E. Darling and Jens Stoye, editors, *Proceedings of the 13th International Workshop on Algorithms in Bioinformatics (WABI '13)*, pages 244–258, 2013.
- 9 Laurent Bulteau and Christian Komusiewicz. Minimum common string partition parameterized by partition size is fixed-parameter tractable. In Chandra Chekuri, editor, *Proceedings of the Twenty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms, (SODA '14)*, pages 102–121. SIAM, 2014.
- 10 Wenbin Chen, Zhengzhang Chen, Nagiza F. Samatova, Lingxi Peng, Jianxiong Wang, and Maobin Tang. Solving the maximum duo-preservation string mapping problem with linear programming. *Theor. Comput. Sci.*, 530:1–11, 2014.
- 11 Xin Chen, Jie Zheng, Zheng Fu, Peng Nan, Yang Zhong, Stefano Lonardi, and Tao Jiang. Assignment of orthologous genes via genome rearrangement. *IEEE/ACM T. Comput. Bi.*, 2(4):302–315, 2005.
- 12 Graham Cormode and S. Muthukrishnan. The string edit distance matching problem with moves. *ACM T. Alg.*, 3(1):2:1–2:19, 2007.
- 13 Maxime Crochemore, Christophe Hancart, and Thierry Lecroq. *Algorithms on strings*. Cambridge University Press, 2007.
- 14 Marek Cygan, Fedor V. Fomin, Lukasz Kowalik, Daniel Lokshtanov, Dániel Marx, Marcin Pilipczuk, Michal Pilipczuk, and Saket Saurabh. *Parameterized Algorithms*. Springer, 2015.
- 15 Peter Damaschke. Minimum common string partition parameterized. In Keith A. Crandall and Jens Lagergren, editors, *Proceedings of the 8th International Workshop on Algorithms in Bioinformatics (WABI '08)*, volume 5251 of *Lecture Notes in Computer Science*, pages 87–98. Springer, 2008.
- 16 Richard A. DeMillo and Richard J. Lipton. A probabilistic remark on algebraic program testing. *Inf. Process. Lett.*, 7(4):193–195, 1978.
- 17 Rodney G. Downey and Michael R. Fellows. *Fundamentals of Parameterized Complexity*. Texts in Computer Science. Springer, 2013.

- 18 Bartłomiej Dudek, Paweł Gawrychowski, and Piotr Ostropolski-Nalewaja. A family of approximation algorithms for the maximum duo-preservation string mapping problem. *CoRR*, arXiv:1702.02405, 2017.
- 19 Guillaume Fertin, Anthony Labarre, Irena Rusu, Eric Tannier, and Stéphane Vialette. *Combinatorics of Genome Rearrangements*. Computational Molecular Biology. MIT Press, 2009.
- 20 Fedor V. Fomin, Daniel Lokshtanov, Fahad Panolan, and Saket Saurabh. Efficient computation of representative families with applications in parameterized and exact algorithms. *J. ACM*, 63(4):29:1–29:60, 2016.
- 21 Avraham Goldstein, Petr Kolman, and Jie Zheng. Minimum common string partition problem: Hardness and approximations. *Electron. J. Comb.*, 12, 2005.
- 22 Dan Gusfield. *Algorithms on Strings, Trees, and Sequences - Computer Science and Computational Biology*. Cambridge University Press, 1997.
- 23 Haitao Jiang, Binhai Zhu, Daming Zhu, and Hong Zhu. Minimum common string partition revisited. *J. Comb. Optim.*, 23:519–527, 2012.
- 24 Jacob T. Schwartz. Fast probabilistic algorithms for verification of polynomial identities. *J. ACM*, 27(4):701–717, 1980.
- 25 Hadas Shachnai and Meirav Zehavi. Representative families: A unified tradeoff-based approach. *J. Comput. Syst. Sci.*, 82(3):488–502, 2016.
- 26 Krister M. Swenson, Mark Marron, Joel V. Earnest-DeYoung, and Bernard M. E. Moret. Approximating the true evolutionary distance between two genomes. *ACM J. Exp. Alg.*, 12, 2008.
- 27 Yao Xu, Yong Chen, Taibo Luo, and Guohui Lin. A local search 2.917-approximation algorithm for duo-preservation string mapping. *CoRR*, arXiv:1702.01877, 2017.
- 28 Richard Zippel. Probabilistic algorithms for sparse polynomials. In Edward W. Ng, editor, *Proceedings of the International Symposium on Symbolic and Algebraic Computation (EUROSAM '79)*, volume 72 of *Lecture Notes in Computer Science*, pages 216–226. Springer, 1979.

A Proofs of Section 3

► **Lemma 14.** *Given $L \subseteq [r]$ and an assignment to all variables $x_{i,j}$ and y_e , the polynomial POL_L can be evaluated in time $O(r^{O(1)} \cdot |E(G)|)$.*

Sketch. The evaluation can be performed by a simple procedure based on dynamic programming. For the sake of completeness, we present the base cases and recursive formula below. For simplicity, we abuse notation by using the symbols $x_{i,j}$ and y_e to refer to the values assigned to the variables $x_{i,j}$ and y_e , respectively.

The procedure uses a table M , which has an entry $M[v, s, b]$ for all $v \in V(G)$, $s \in [r]$ and $b \in [r] \cup \{0\}$. The purpose of this entry is to store the evaluation of $\text{POL}_{v,s,b,L}$. Then, the evaluation of POL_L is given by

$$\sum_{\substack{v \in V(G), \\ s, b \in [r], b \geq k}} M[v, s, b].$$

The basis consists of the following cases:

- If $b \geq s$, then $M[v, s, b] = 0$.
- Else if $s = 1$, then $M[v, s, b] = \sum_{i \in L} x_{\ell(v), i}$.

Now, consider an entry $M[v, s, b]$ not computed in the basis. We assume that a reference to an undefined entry returns 0. Then,

$$\begin{aligned} M[v, s, b] = & \sum_{\substack{(u, v) \in E(G), \\ c(u, v) = R}} \left(\sum_{i \in L} x_{\ell(v), i} \cdot y_{(u, v)} \cdot M[u, s - 1, b] \right) \\ & + \sum_{\substack{(u, v) \in E(G), \\ c(u, v) = B}} \left(\sum_{i \in L} x_{\ell(v), i} \cdot y_{(u, v)} \cdot M[u, s - 1, b - 1] \right). \end{aligned}$$

◀

B Proofs of Section 4

In this section, we adapt the approach in which the method of *representative sets* is applied to solve the k -PATH problem [20]. More precisely, our objective is to provide a constructive proof for the following result.

► **Lemma 15.** *There exists a deterministic algorithm that solves LONG BLUE PATH in time $O\left(\left(\frac{1+\sqrt{5}}{2}\right)^{r+o(r)} \cdot |E(G)| \log |E(G)|\right)$.*

In light of Lemma 3, once we have Lemma 15 at hand, we directly obtain the following theorem.

► **Theorem 16.** *There exists a deterministic algorithm that solves MAX-DUO PSM in time $O\left(\left(\frac{1+\sqrt{5}}{2}\right)^{2k+o(k)} \cdot n^4 \log n\right) = O(6.855^k \cdot n^4 \log n)$.*

Next, we focus on the proof of Lemma 15. To this end, let (G, c, ℓ, k, r) be an instance of LONG BLUE PATH. Without loss of generality, we can assume that the image of ℓ is a subset of $[V(G)]$ and that $|V(G)| \leq |E(G)|$. Here, a p -set is a set of size p . To describe our algorithm, we need to present the definition of a representative family.

► **Definition 27** ([20]). Given a universe U and a family \mathcal{S} of p -subsets of U , we say that a subfamily $\widehat{\mathcal{S}} \subseteq \mathcal{S}$ t -represents \mathcal{S} if for every pair of sets $X \in \mathcal{S}$, and $Y \subseteq U \setminus X$ of size $t - p$, there exists a set $\widehat{X} \in \widehat{\mathcal{S}}$ such that $\widehat{X} \cap Y = \emptyset$.

The papers [20] and [25] present an algorithm, to which we refer as **RepAlg**, that given a universe U and a family \mathcal{S} of p -subsets of U , computes a subfamily $\widehat{\mathcal{S}} \subseteq \mathcal{S}$ of size $S(|U|, t, p)$ that t -represents \mathcal{S} in time $|\mathcal{S}| \cdot T(|U|, t, p)$, such that the following condition is satisfied:

$$\sum_{p=1}^t |U| \cdot S(|U|, t, p - 1) \cdot T(|U|, t, p) = \left(\frac{1 + \sqrt{5}}{2}\right)^{t+o(t)} \cdot |U| \log |U|.$$

We proceed by presenting a procedure that is based on a combination of dynamic programming and calls to **RepAlg**. For this purpose, we use a table M that has an entry $M[v, s, b]$ for all $v \in V(G)$, $s \in [r]$ and $b \in [r] \cup \{0\}$. Let $\mathcal{P}_{v,s,b}$ denote the set of all *good* (v, s, b) -paths (see Definition 6). Give a (v, s, b) -path, define $\ell(P) = \{\ell(v) : v \in V(P)\}$. Moreover, define $\mathcal{S}_{v,s,b} = \{\ell(P) : P \in \mathcal{P}_{v,s,b}\}$. The purpose of the entry $M[v, s, b]$ would be to store a subfamily of $\mathcal{S}_{v,s,b}$ that r -represents it. Next, we show how to compute the entries of M . Here, the calls to **RepAlg** correspond to the universe $[|E(G)|]$ and with $t = r$.

The basis consists of the following cases:

- If $s = 1$ but $b \neq 0$, then $M[v, s, b] = \emptyset$.
- Else if $s = 1$, then $M[v, s, b] = \{\{\ell(v)\}\}$.

Now, consider an entry $M[v, s, b]$ not computed in the basis. We assume that a reference to an undefined entry returns an empty set. Then, we first compute the two following families.

- $\mathcal{A}_{v,s,b} = \{X \cup \{\ell(v)\} : (u, v) \in E(G), c(u, v) = R, X \in M[u, s - 1, b], \ell(v) \notin X\}$.
- $\mathcal{B}_{v,s,b} = \{X \cup \{\ell(v)\} : (u, v) \in E(G), c(u, v) = B, X \in M[u, s - 1, b - 1], \ell(v) \notin X\}$.

Accordingly, we compute $M[v, s, b]$ as follows.

$$M[v, s, b] = \text{RepAlg}(\mathcal{A}_{v,s,b} \cup \mathcal{B}_{v,s,b}).$$

First, note that the entire computation can be performed in time

$$\begin{aligned} & O\left(\sum_{v \in V(G)} \sum_{s=1}^r \sum_{b=0}^r \sum_{(u,v) \in E(G)} S(|E(G)|, r, s) \cdot T(|E(G)|, r, s)\right) \\ &= O\left(\sum_{s=1}^r r |E(G)| \cdot S(|E(G)|, r, s) \cdot T(|E(G)|, r, s)\right). \end{aligned}$$

Thus, we have the following observation.

► **Observation 7.** *The table M is computed in time $O\left(\left(\frac{1+\sqrt{5}}{2}\right)^{r+o(r)} \cdot |E(G)| \log |E(G)|\right)$.*

Next, we prove that the computation of M is correct.

► **Lemma 28.** *The computation of M ensures that for all $v \in V(G)$, $s \in [r]$ and $b \in [r] \cup \{0\}$, $M[v, s, b]$ r -represents $\mathcal{S}_{v,s,b}$.*

Proof. We prove the statement by induction on s . In the basis, where $s = 1$, it is clear that $M[v, s, b]$ is simply assigned $\mathcal{S}_{v,s,b}$, and therefore it also 1-represents $\mathcal{S}_{v,s,b}$. Now, fix some $s \geq 2$, and suppose that the statement is correct for $s - 1$. To prove that the statement is correct for s , choose some $v \in V(G)$, $b \in [r] \cup \{0\}$, $X \in \mathcal{S}_{v,s,b}$ and $Y \subseteq [|E(G)|] \setminus X$ such that

$|Y| = r - s$. We need to show that there exists $\widehat{X} \in M[v, s, b]$ such that $\widehat{X} \cap Y = \emptyset$. Note that $M[v, s, b]$ r -represents $\mathcal{A}_{v,s,b} \cup \mathcal{B}_{v,s,b}$, and therefore $\mathcal{A}_{v,s,b} \cup \mathcal{B}_{v,s,b}$ contains a set that is disjoint from Y , so does $M[v, s, b]$. Thus, it is sufficient that we show that there exists $\widehat{X} \in \mathcal{A}_{v,s,b} \cup \mathcal{B}_{v,s,b}$ such that $\widehat{X} \cap Y = \emptyset$.

Since $X \in \mathcal{S}_{v,s,b}$, there exists a good (v, s, b) -path P such that $\ell(P) = X$. Let u be the vertex on P that precedes v , and let Q be the path obtained by removing v from P . Note that $\ell(Q) = X \setminus \{\ell(v)\}$. Thus, if $c(u, v) = R$, then Q is a good $(u, s - 1, b)$ -path and therefore $X \setminus \{\ell(v)\} \in \mathcal{S}_{u,s-1,b}$, and otherwise Q is a good $(u, s - 1, b - 1)$ -path and therefore $X \setminus \{\ell(v)\} \in \mathcal{S}_{u,s-1,b-1}$. First, let us assume that $X \setminus \{\ell(v)\} \in \mathcal{S}_{u,s-1,b}$. By the inductive hypothesis, $M[u, s - 1, b]$ r -represents $\mathcal{S}_{u,s-1,b}$, and therefore $M[u, s - 1, b]$ contains a set Z such that $Z \cap (Y \cup \{\ell(v)\}) = \emptyset$. Thus, $Z \cup \{\ell(v)\} \in \mathcal{A}_{v,s,b}$, and we conclude that the statement is correct. Now, let us assume that $X \setminus \{\ell(v)\} \in \mathcal{S}_{u,s-1,b-1}$. By the inductive hypothesis, $M[u, s - 1, b - 1]$ r -represents $\mathcal{S}_{u,s-1,b-1}$, and therefore $M[u, s - 1, b - 1]$ contains a set Z such that $Z \cap (Y \cup \{\ell(v)\}) = \emptyset$. Thus, $Z \cup \{\ell(v)\} \in \mathcal{B}_{v,s,b}$, and again we conclude that the statement is correct. ◀

With these lemmas at hand, we are ready to prove Lemma 15.

Proof. By Observation 7 and Lemma 28, we first compute M , ensuring that the condition in Lemma 28 is satisfied, in time $O((\frac{1+\sqrt{5}}{2})^{r+o(r)} \cdot |E(G)| \log |E(G)|)$. Then, we determine that the input instance is a yes-instance if and only if there exist $v \in V(G)$, $s \in [r]$ and $b \in [r]$ such that $b \geq k$ and $M[v, s, b] \neq \emptyset$. On the one hand, since for all $v \in V(G)$, $s \in [r]$ and $b \in [r]$, $M[v, s, b] \subseteq \mathcal{S}_{v,s,b}$, it is clear that if we accept, the input instance is indeed a yes-instance. On the other hand, if the input instance is a yes-instance, then there exist $v \in V(G)$, $s \in [r]$ and $b \in [r]$ such that $b \geq k$ and $\mathcal{S}_{v,s,b} \neq \emptyset$. Then, since $M[v, s, b]$ 0-represents $\mathcal{S}_{v,s,b}$, it holds that $M[v, s, b] \neq \emptyset$, and therefore we accept. ◀

C Proofs of Section 5

Proof of Theorem 23 We first show the running time to construct the kernel.

▶ **Proposition 1.** *For each $S \in \{A, B\}$, $|C_S| = O(k^3)$ and C_S can be constructed in time $O(|\Sigma|^2 \cdot n) + O(k^3 n)$.*

Proof. By assumption $|rare(S)| \leq 4k$. Additionally, for each i , the ball $\mathcal{B}_k(i)$ has size at most $2k + 1$. Finally, for each duo $(i, i + 1)$ that is rare for S , the set $F(\overline{S}, i, i + 1)$ has at most $(2k - 1)(2k + 1)$ duos. Therefore, $|C_S| \leq 4k(2k + 1) + 4k(2k - 1)(2k + 1) = O(k^3)$.

Now let us analyze the time to construct C_S . First, the construction of the sets $rare(S)$ and $rare(\overline{S})$ takes time $O(|\Sigma|^2 \cdot n)$, since we just need to count for each length-two string $ab \in \Sigma \times \Sigma$, the number of times $n(S, a, b)$ that ab occurs in S and the number of times $n(\overline{S}, a, b)$ that ab occurs in \overline{S} . Now, for each position $i \in \{1, \dots, n - 1\}$, we add $(i, i + 1)$ to $rare(S)$ if $S[i]S[i + 1] = ab$ and $n(S, a, b) \leq n(\overline{S}, a, b)$. Analogously, we add $(i, i + 1)$ to $rare(\overline{S})$ if $\overline{S}[i]\overline{S}[i + 1] = ab$ and $n(\overline{S}, a, b) \leq n(S, a, b)$.

Now, the construction of the set $ROOTS(S, i, i + 1)$ according to Algorithm 1 takes time $O(k^2 \cdot n)$. Since $ROOTS(S, i, i + 1) \leq 2k - 1$, and by assumption $|rare(S)| \leq 4k$, the construction of $F(S, i, i + 1)$ also takes time $O(k^2 \cdot n)$. Analogously, the construction of $F(\overline{S}, i, i + 1)$ takes time $O(k^2 \cdot n)$. Therefore, the construction of C_S takes time at most $O(|\Sigma|^2 \cdot n) + O(k^3 \cdot n)$. ◀

► **Theorem 23** . *Given an instance $I = (A, B, k)$ of MAX-DUO PSM, one can construct in time $O(|\Sigma|^2 \cdot n + k^3 \cdot n)$ an instance $I'(A', B', k)$ of MAX-DUO PSM with $|A'|$ and $|B'|$ bounded by $O(k^3)$ such that I is a yes-instance if and only if I' is a yes-instance.*

Proof. First, if some map m of A into B preserves a block X of size k , then (A, B, k) is a yes-instance for MAX-DUO PSM and we can output in $O(1)$ time an equivalent instance of constant size. Note that this condition can be verified in time $O(n)$ by solving the LONGEST COMMON SUBSTRING problem for A and B .

Second, if $\text{rare}(A) \geq 4k$ or $\text{rare}(B) \geq 4k$, then (A, B, k) is a yes-instance for MAX-DUO PSM, and we can output in $O(1)$ time an equivalent instance of constant size.

Now since no map preserves a block of size k , and if both $\text{rare}(A) < 4k$ and $\text{rare}(B) < 4k$, then by Lemma 22, the pair (C_A, C_B) constructed according to Equation 1 is complete for (A, B, k) . Additionally, by Proposition 1, $|C_A| = |C_B| = O(k^3)$, and both C_A and C_B can be constructed in time $O(|\Sigma| \cdot n + k^3 \cdot n)$.

Since the complete pair (C_A, C_B) constructed has size at most $O(k^3)$, we can apply Theorem 17 to construct in time $O(k^3)$ an instance (A', B', k) for MAX-DUO PSM of size $O(k^3)$ such that (A', B', k) is a yes-instance if and only if (A, B, k) is a yes-instance. Therefore, the overall time to construct (A', B', k) is upper-bounded by $O(|\Sigma|^2 \cdot n + k^3 \cdot n)$. ◀