

Edit Distance Cannot Be Computed in Strongly Subquadratic Time (unless SETH is false)

Arturs Backurs*
MIT

Piotr Indyk†
MIT

Abstract

The edit distance (a.k.a. the Levenshtein distance) between two strings is defined as the minimum number of insertions, deletions or substitutions of symbols needed to transform one string into another. The problem of computing the edit distance between two strings is a classical computational task, with a well-known algorithm based on dynamic programming. Unfortunately, all known algorithms for this problem run in nearly quadratic time.

In this paper we provide evidence that the near-quadratic running time bounds known for the problem of computing edit distance might be tight. Specifically, we show that, if the edit distance can be computed in time $O(n^{2-\delta})$ for some constant $\delta > 0$, then the satisfiability of conjunctive normal form formulas with N variables and M clauses can be solved in time $M^{O(1)}2^{(1-\epsilon)N}$ for a constant $\epsilon > 0$. The latter result would violate the *Strong Exponential Time Hypothesis*, which postulates that such algorithms do not exist.

*backurs@mit.edu

†indyk@mit.edu

1 Introduction

The edit distance (a.k.a. the Levenshtein distance) between two strings is defined as the minimum number of insertions, deletions or substitutions of symbols needed to transform one string into another. It is a widely used metric, with many applications in computational biology, natural language processing and information theory. The problem of computing the edit distance between two strings is a classical computational task, with a well-known algorithm based on dynamic programming. Unfortunately, that algorithm runs in quadratic time, which is prohibitive for long sequences (e.g., the human genome consists of roughly 3 billions base pairs). A considerable effort has been invested into designing faster algorithms, either by assuming that the edit distance is bounded, by considering the average case or by resorting to approximation¹. However, the fastest known exact algorithm, due to [MP80], has a running time of $O(n^2/\log^2 n)$ for sequences of length n , which is still nearly quadratic.

In this paper we provide evidence that the (near)-quadratic running time bounds known for this problem might, in fact, be tight. Specifically, we show that if the edit distance can be computed in time $O(n^{2-\delta})$ for some constant $\delta > 0$, then the satisfiability of conjunctive normal form (CNF) formulas with N variables and M clauses can be solved in time $M^{O(1)}2^{(1-\epsilon)N}$ for a constant $\epsilon > 0$. The latter result would violate the *Strong Exponential Time Hypothesis (SETH)*, introduced in [IP01, IPZ01], which postulates that such algorithms do not exist². The rationale behind this hypothesis is that, despite decades of research on fast algorithms for satisfiability and related problems, no algorithm was yet shown to run in time faster than $2^{N(1-o(1))}$. Because of this state of affairs, SETH has served as the basis for proving conditional lower bounds for several important computational problems, including k-Dominating Set [PW10], the diameter of sparse graphs [RW13], local alignment [AWW14], dynamic connectivity problems [AW14], and the Frechet distance computation [Bri14]. Our paper builds on these works, identifying a new important member of the class of “SETH-hard” problems.

Our techniques and related work This work has been stimulated by the recent result of [Bri14], who showed an analogous hardness result for computing the Frechet distance³, and listed SETH-hardness of edit distance as an open problem. There are notable similarities between the edit distance and the Frechet distance. In particular, both can be computed in quadratic time, via dynamic programming over an $n \times n$ table T where each entry $T[i, j]$ holds the distance between the first i elements of the first sequence and the first j elements of the second sequence. Furthermore, in both cases each entry $T[i, j]$ can be computed locally given $T[i, j-1]$, $T[i-1, j]$ and $T[i-1, j-1]$. The key difference between the two distances is that while the recursive formula for the Frechet distance uses the max function, the formula for the edit distance involves the sum. As a result, the Frechet distance is effectively determined by a single pair of sequence elements, while the edit distance is determined by many pairs of elements. As we describe below, this makes the reduction

¹There is a rather large body of work devoted to edit distance algorithms and we will not attempt to list all relevant works here. Instead, we refer the reader to the survey [Nav01] for a detailed overview of known exact and probabilistic algorithms, and to the recent paper [AKO10] for an overview of approximation algorithms.

²Technically, our results relies on an even weaker conjecture. See Preliminaries for more details.

³Given two sequences of points P_1 and P_2 , the Frechet distance between them is defined as the minimum, over all monotone traversals of P_1 and P_2 , of the largest distance between the corresponding points at any stage of the traversal.

to edit distance much more subtle.⁴

Our result is obtained by a reduction from the *Orthogonal Vectors Problem (OVP)*, which is defined as follows. Given two sets $A, B \subseteq \{0, 1\}^d$ such that $|A| = |B| = N$, the goal is to determine whether there exists $x \in A$ and $y \in B$ such that the dot product $x \cdot y = \sum_{j=1}^d x_j y_j$ (taken over reals) is equal to 0. It is known [Wil05] that an $O(N^{2-\delta} d^{O(1)})$ -time algorithm for OVP would imply that SETH is false. Therefore, in what follows we focus on reducing Orthogonal Vectors problem to the Edit Distance problem.

The first step of our reduction mimics the approaches in [Bri14] and [AWW14]. In particular, each $x \in A$ and $y \in B$ is assigned a “gadget” sequence. Then, the gadget sequences for all $a \in A$ are concatenated together to form the first input sequence, and the gadget sequences for all $b \in B$ are concatenated to form the second input sequence. The correctness of the reduction is proven by showing that:

- If there is a pair of orthogonal vectors $x \in A$ and $y \in B$, then one can traverse the two sequences in a way that the gadgets assigned to x and y are *aligned*, which implies that the distance induced by this traversal is “small”.
- If there is no orthogonal pair, then no such traversal exists, which implies that the distance induced by any traversal is “large”.

The mechanics of this argument depends on the specific distance function. In the case of Frechet distance, the output value is determined by the maximum distance between the aligned elements, so it suffices to show that the distance between two vector gadgets is smaller than C if they are orthogonal and at least C if they are not, for some value of C . In contrast, edit distance *sums up* the distances between *all* aligned gadgets (as well as the costs of insertions and deletions used to create the alignment), which imposes stronger requirements on the construction. Specifically, we need to show that if two vectors x and y are not orthogonal, i.e., they have at least one overlapping 1, then the distance between their gadgets is *equal* to C , not just at least C . Since we need to ensure that the distance between two gadgets cannot grow in the number of overlapping 1s, our gadget design and analysis become more complex.

Fortunately, the edit distance is expressive enough to support this functionality. The basic idea behind the gadget construction is to use that fact the edit distance between two gadget strings, say AG_1 (from the first sequence) and AG_2 (from the second sequence), is the minimum cost over all possible alignments between AG_1 and AG_2 . Specifically, we construct gadgets that allow two alignment options. The first option results in a cost that is linear in the number of overlapping 1s of the corresponding vectors (this is easily achieved by using substitutions only). On the other hand, the second “fallback” option has a fixed cost (say C_1) that is slightly higher than the cost of the first option when no 1s are overlapping (say, C_0). Thus, by taking the minimum of these two options, the resulting cost is equal to C_0 when the vectors are orthogonal and equal to C_1 ($> C_0$) otherwise, which is what is needed. See Lemmas 1 and 2 for the details of the construction.

Further developments Following this work, two recent technical reports showed multiple results demonstrating conditional hardness of the edit distance, the longest common subsequence problem, dynamic time warping problem and other similarity measures between sequences [ABW15, BK15].

⁴This also means that our hardness argument does not extend to the approximate edit distance computation, in contrast to the argument in [Bri14].

Among other results, [BK15] showed hardness of computing the edit distance over the *binary* alphabet, which improves over the alphabet size of 4 required for our reduction.

2 Preliminaries

Edit distance For any two sequences x and y over an alphabet Σ , the edit distance $\text{EDIT}(x, y)$ is equal to the minimum number of symbol insertions, symbol deletions or symbol substitutions needed to transform x into y . It is well known that the EDIT function induces a metric; in particular, it is symmetric and satisfies the triangle inequality.

In the remainder of this paper we will use an equivalent definition of EDIT that will make the analysis of our reductions more convenient.

Observation 1. *For any two sequences x, y , $\text{EDIT}(x, y)$ is equal to the minimum, over all sequences z , of the number of deletions and substitutions needed to transform x into z and y into z .*

Proof. It follows directly from the metric properties of the edit distance that $\text{EDIT}(x, y)$ is equal to the minimum, over all sequences z , of the number of *insertions*, deletions and substitutions needed to transform x into z and y into z . Furthermore, observe that if, while transforming x , we insert a symbol that is later aligned to some symbol of y , we can instead delete the corresponding symbol in y . Thus, it suffices to allow deletions and substitutions only. \square

Definition 1. *We define the following similarity distance between sequences P_1 and P_2 and we call it the pattern matching distance between P_1 and P_2 .*

$$\text{PATTERN}(P_1, P_2) = \min_{\substack{x \text{ is a contiguous} \\ \text{subsequence of } P_2}} \text{EDIT}(P_1, x).$$

For a symbol a and an integer i we use a^i to denote symbol a repeated i times.

Orthogonal vectors problem The *Orthogonal Vectors Problem* is defined as follows: given two sets $A, B \subseteq \{0, 1\}^d$ such that $|A| = |B| = N$, determine whether there exists $x \in A$ and $y \in B$ such that the dot product $x \cdot y = \sum_{j=1}^d x_j y_j$ (taken over reals) is equal to 0. An equivalent formulation of this problem is: given two collections of n sets each, determine if there a set in the first collection that is contained in a set from the second collection.

The orthogonal vectors problem has an easy $O(N^2 d)$ -time solution. However, it is known that any algorithm for this problem with strongly sub-quadratic running time would also yield a more efficient algorithm for SAT, breaking SETH [Wil05]. Thus, in what follows, we focus on reducing Orthogonal Vectors Problem to EDIT.

3 Reductions

We will first show that solving Orthogonal Vectors Problem can be reduced to computation of PATTERN. Then we will show how the computation of PATTERN can be reduced to computation of EDIT.

3.1 Reducing Orthogonal Vectors Problem to PATTERN

We now describe vector gadgets as well as provide an intuition behind the construction.

We will construct sequences over an alphabet $\Sigma = \{0, 1, 2\}$.

We start by defining integer parameters $l_0 = 1000 \cdot d$ and $l_1 = (1000 \cdot d)^2$. We then define *coordinate gadget* sequences CG_1 and CG_2 as follows. For integer $x \in \{0, 1\}$ we define

$$\text{CG}_1(x) := \begin{cases} 0^{l_1} 0^{l_0} 1^{l_0} 1^{l_0} 1^{l_0} 0^{l_1} & \text{if } x = 0 \\ 0^{l_1} 0^{l_0} 0^{l_0} 0^{l_0} 1^{l_0} 0^{l_1} & \text{if } x = 1 \end{cases}$$

$$\text{CG}_2(x) := \begin{cases} 0^{l_1} 0^{l_0} 0^{l_0} 1^{l_0} 1^{l_0} 0^{l_1} & \text{if } x = 0 \\ 0^{l_1} 1^{l_0} 1^{l_0} 1^{l_0} 1^{l_0} 0^{l_1} & \text{if } x = 1 \end{cases}$$

We will also need the following sequence:

$$g = 0^{\frac{l_1}{2}-1} 1 0^{\frac{l_1}{2}} 0^{l_0} 1^{l_0} 1^{l_0} 1^{l_0} 0^{l_1}.$$

The distance properties of the defined sequences are depicted in Figure 1.

The coordinate gadgets were designed so that they have the following properties. For any two integers $x_1, x_2 \in \{0, 1\}$,

$$\text{EDIT}(\text{CG}_1(x_1), \text{CG}_2(x_2)) = \begin{cases} l_0 & \text{if } x_1 \cdot x_2 = 0 \\ 3l_0 & \text{if } x_1 \cdot x_2 = 1. \end{cases}$$

Also, we have

$$\text{EDIT}(\text{CG}_2(x), g) = l_0 + 1$$

for any integer $x \in \{0, 1\}$.

Further, we define another parameter $l_2 = (1000 \cdot d)^3$. For two sequences x and y , we use $x \circ y$ or $x \circ y$ to denote the concatenation of x and y . For vectors a and b , we define the *vector gadget* sequences as

$$\text{VG}_1(a) = Z_1 L V_0 R Z_2 \text{ and } \text{VG}_2(b) = V_1 D V_2,$$

where

$$Z_1 = Z_2 = 0^{l_2}, \quad V_1 = V_2 = V_0 = 1^{l_2},$$

$$L = \bigcirc_{i \in [d]} g, \quad R = \bigcirc_{i \in [d]} \text{CG}_1(a_i), \quad D = \bigcirc_{i \in [d]} \text{CG}_2(b_i).$$

We denote the length of L , R and D by $l = |L| = |R| = |D| = d(4l_0 + 2l_1)$.

We visualize the defined vector gadgets in Figure 2.

Intuition behind the construction Before going into the analysis of the gadgets in Section 3.1.1, we will first provide some intuition behind the construction. Given two vectors $a, b \in \{0, 1\}^d$, we want that $\text{EDIT}(\text{VG}_1(a), \text{VG}_2(b))$ is small (equal to C_0) if the vectors a and b are orthogonal, and is large (equal to C_1) otherwise. That is, we want that

$$\text{EDIT}(\text{VG}_1(a), \text{VG}_2(b)) = \begin{cases} C_0 & \text{if } a \cdot b = 0 \\ C_1 & \text{otherwise} \end{cases} \quad (1)$$

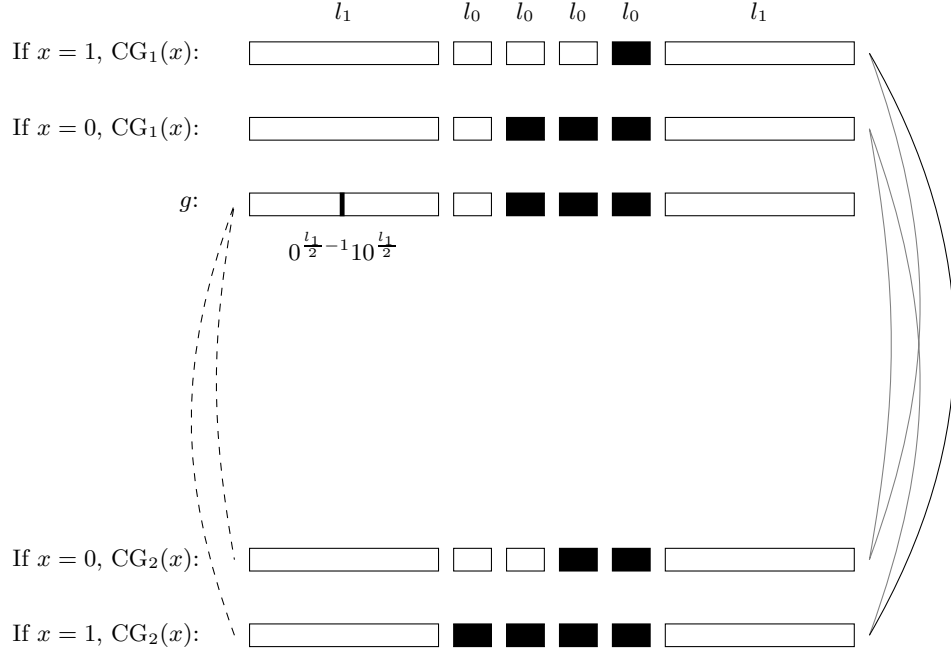


Figure 1: A visualisation of the coordinate gadgets. Each rectangle denotes a run of the same symbol (either 0 or 1). A black rectangle denotes a run of 1s, a white rectangle denotes a run of 0s. A short rectangle denotes a run of length l_0 , while a long rectangle denotes a run of length l_1 . An exception is the first run of g of length l_1 which contains a single 1 in the middle of 0s. A solid black line denotes that the edit distance between the two sequences is equal to $3l_0$; a gray line denotes the edit distance of l_0 and a dashed black line denotes the edit distance of $l_0 + 1$.

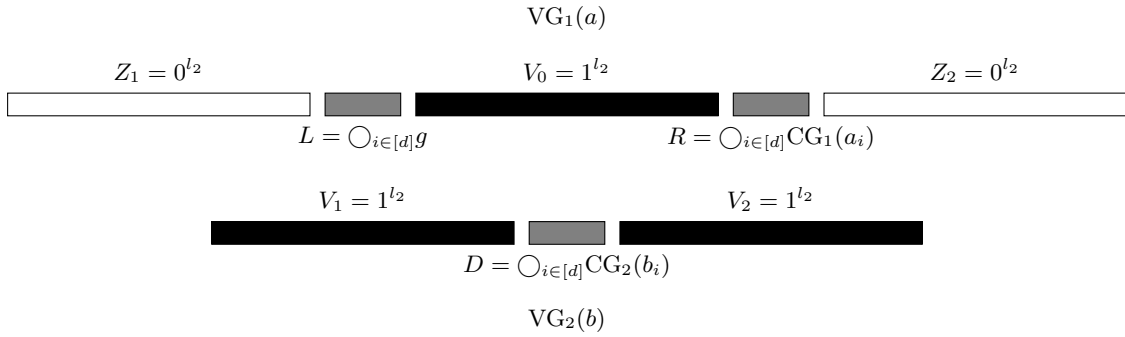


Figure 2: A visualisation of the vector gadgets. A black rectangle denotes a run of 1s, while a white rectangle denotes a run of 0s. A gray rectangle denotes a sequence that contains both 0s and 1s. A short rectangle denotes a sequence of length l , while a long one denotes a sequence of length l_2 .

for $C_1 > C_0$. We also want the property that C_0 and C_1 do not depend on the vectors a and b , that is, C_0 and C_1 are only functions of d . We note that it is easy to construct gadgets with the property that $C_1 = t_0 + t \cdot (a \cdot b)$ for some $t = t(d)$ and $t_0 = t_0(d)$, i.e., such that C_1 is linear in the inner product between the vectors. Unfortunately, this construction does not yield, and cannot possibly yield, the desired reduction (see the note at the end of Section 3.2). In what follows we first describe the construction that achieves $C_1 = t_0 + t \cdot (a \cdot b)$ and then we will show how to fix it.

Our first construction uses $\bar{V}G_1(a) = \bigcirc_{i \in [d]} CG_1(a_i)$ and $\bar{V}G_2(b) = \bigcirc_{i \in [d]} CG_2(b_i)$. It is easy to see that the edit distance between two coordinates is small (l_0) if the vectors are orthogonal in that coordinate and large ($3l_0$) otherwise. Therefore, we can check that $\text{EDIT}(\bar{V}G_1(a), \bar{V}G_2(b)) = t_0 + t \cdot (a \cdot b)$ for some $t_0 = t_0(d)$ and $t = t(d)$.

In order to fix this construction we proceed as follows. Since the edit distance between gadgets $\bar{V}G_1$ and $\bar{V}G_2$ grows linearly in the inner product between the vectors, we need to somehow “cap” this growth. To this end we construct our vector gadgets VG_1 and VG_2 such that there are only two possibilities to achieve small edit distance. In the first case, the edit distance grows linearly in the inner product between the sequences. In the second case, the edit distance is constant, independent from the vectors; in particular it does not depend on whether the vectors are orthogonal or not. Because the edit distance is equal to the minimum over all possible alignments, we take the minimum of a function that grows linearly in the inner product between the vectors and a constant function. After taking the minimum, the edit distance will satisfy the properties stated in (1). More precisely, we achieve the minimum edit distance cost between VG_1 and VG_2 by following one of the following two possible sequences of operations:

- Case 1: Delete Z_1 and L . Substitute Z_2 with V_2 . Transform R and D into the same sequence by transforming corresponding coordinate gadgets into the same sequences. Notice that $R = \bar{V}G_1(a)$ and $D = \bar{V}G_2(b)$. Therefore, this case corresponds to a linearly growing function. The cost of all of these operations is $C_0 + t \cdot (a \cdot b)$ for some constant t .
- Case 2: Delete R and Z_2 . Substitute Z_1 with V_1 . Transform L and D into the same sequence by transforming the corresponding coordinate gadgets. We can check that the value of $\text{EDIT}(L, D)$ does not depend on the input vectors. Thus, this case corresponds to achieving constant edit distance cost. The cost of all of these operations is $C_1 (> C_0)$.

Taking the minimum of these two cases yields the desired formula (1). The sequences Z_1, Z_2, V_0, V_1, V_2 are needed to implement the minimum operation. In the next section we give the details of the described construction. In fact, we show a somewhat weaker property: if $a \cdot b = 0$, then $\text{EDIT}(VG_1(a), VG_2(b)) \leq C_0$. However, this change does not affect the validity of the reduction.

3.1.1 Gadget properties

For a binary sequence x , we use $\text{ones}(x)$ ($\text{zeros}(x)$) to denote number of 1s (0s) in x . The vector gadgets have the following important properties.

If two vectors are orthogonal, then the edit distance between them is E_s , which is only a function of d .

Lemma 1. *If $a \cdot b = 0$ (that is, vectors a and b are orthogonal), then*

$$\text{EDIT}(VG_1(a), VG_2(b)) \leq 2l_2 + l + dl_0 =: E_s.$$

Proof. We can transform $VG_1(a)$ and $VG_2(b)$ into the same sequence as follows. First, we delete Z_1 and L from $VG_1(a)$, incurring a cost of $l_2 + l$. Second, the remaining sequence V_0RZ_2 can be transformed into $VG_2(b)$ using $l_2 + dl_0$ substitutions. \square

If two vectors are not orthogonal, then the edit distance between them is E_u ($> E_s$), which is only a function of d .

Lemma 2. *If $a \cdot b \geq 1$ (that is, vectors a and b are not orthogonal), then*

$$EDIT(VG_1(a), VG_2(b)) = l + 2l_2 + dl_0 + d =: E_u.$$

Proof. Observe that $EDIT(VG_1(a), VG_2(b)) \leq E_u$. To achieve this bound we first delete R and Z_2 from $VG_1(a)$, which costs $l_2 + l$. The remaining sequence Z_1LV_0 can be transformed into $VG_2(b)$ using $l_2 + dl_0 + d$ substitutions.

In what follows, we show that $EDIT(VG_1(a), VG_2(b)) \geq E_u$. Consider an optimal transformation of $VG_1(a)$ and $VG_2(b)$ into the same sequence. Every symbol in the first sequence is either substituted, preserved or deleted in the process. If a symbol is not deleted, we call it an *alignment* between the corresponding pair of symbols from the sequences.

We state a simple fact without proof that we will use multiple times later.

Fact 1. *Suppose we have two sequences x and y of symbols. Let $i_1 < j_1$ and $i_2 < j_2$ be four positive integers. If x_{i_1} is aligned with y_{j_2} , then x_{j_1} can't be aligned with y_{i_2} .*

From now on we proceed by considering three main cases.

Case 1. The subsequence D has alignments with both Z_1L and RZ_2 . In this case, the cost induced by symbols from Z_1 and Z_2 , and V_0 is $\geq l_2 - l$ for each one of these sequences. This implies that $EDIT(VG_1(a), VG_2(b)) \geq 3l_2 - l$, which contradicts an easy upper bound. We have an upper bound $EDIT(VG_1(a), VG_2(b)) \leq 2l_2 + 3l$, which is obtained by deleting L, R, D, Z_1 and replacing Z_2 with V_2 symbol by symbol. Therefore, this case cannot occur.

Case 2. D does not have any alignments with Z_1L . In this case we start by introducing the following notion. Let v and z be two sequences that decompose as $v = xV$ and $z = yZ$. Consider two sequences T (R , resp.) of deletions and substitutions that transform v into u (z into u , resp.). An operation in T or R is called *internal to V and Z* if it is either a (1) deletion of a symbol in V or Z , or (2) a substitution of a symbol in V so that it aligns with a symbol in Z , or vice versa. All other operations, including substitutions that align with symbols in V (Z , resp.) to those outside of Z (V , resp.) are called *external to V and Z* .

We state the following simple fact without proof.

Fact 2. *Let xV and yZ be sequences such that $|V| = |Z| = t$, V consists of only 1s, Z consists of only 0s, and x and y are arbitrary sequences over an arbitrary alphabet. Consider $EDIT(xV, yZ)$ and the corresponding operations minimizing the distance. Among those operations, the number of operations that are internal to V and Z is at least t .*

Given that $|Z_2| = |V_2| = l_2$ and Z_2 consists of only 0s and V_2 consists of only 1s, Fact 2 implies that the number of operations that are internal to Z_2 and V_2 is at least $S_1 := l_2$. In order to lower bound the number of operations that are external to Z_2 and V_2 , we consider the following four sub-cases. First, we consider the case when there is a symbol 1 in V_0 that is aligned with a symbol 1 in DV_2 (Case 2.1). Second, we consider the case when there is a symbol 1 in V_1 that is aligned with symbol 1 in R (Case 2.2). Third, there is a symbol 1 in V_2 that is aligned with a symbol 1

in R (Case 2.3). Finally, we consider the case when, if a symbol 1 in R (D , resp.) is aligned with a symbol 1, then the latter symbol comes from D (R , resp.) (Case 2.4). We can verify that these cases exhaust all possible scenarios. We proceed by analyzing the listed cases in detail.

Case 2.1. A symbol 1 in V_0 is aligned with a symbol 1 in DV_2 . This implies that the first l_1 symbols 0 of D are deleted or flipped (we use the fact that D does not have alignments with Z_1L). The cost of this is at least l_1 , which is more than $S_2 := \text{ones}(L) + dl_0 + d$. Note also that there are no 0s in $VG_2(b)$ to be aligned with 0s that are to the left of V_0 . Because of this, each 0 to the left of V_0 induces a cost of 1 per symbol. The number of 0s to the left of V_0 is $S_3 := \text{zeros}(Z_1L)$, which is also the lower bound for the induced cost. Summing up all the costs, we obtain that the $\text{EDIT}(VG(a), VG(b))$ is at least

$$S_1 + S_2 + S_3 = l_2 + (\text{ones}(L) + dl_0 + d) + \text{zeros}(Z_1L) = E_u.$$

Case 2.2. A symbol 1 in V_1 is aligned with symbol 1 in R . This implies that the first $l_1 \geq \text{ones}(L) + dl_0 + d = S_2$ symbols 0 of R are deleted or flipped. Because there are no 0s in $VG_2(b)$ to be aligned with 0s that are to the left of V_0 , every zero that is to the left of V_0 induces a cost of at least 1. The total induced cost is at least $S_3 = \text{zeros}(Z_1L)$. We conclude that

$$\begin{aligned} \text{EDIT}(VG(a), VG(b)) &\geq S_1 + S_2 + S_3 \\ &= l_2 + (\text{ones}(L) + dl_0 + d) + \text{zeros}(Z_1L) = E_u. \end{aligned}$$

Therefore, from now on we assume that neither Case 2.1 or Case 2.2 occur, i.e., no symbol of V_0 or V_1 is aligned with a symbol 1 in DV_2 or R , respectively. From this assumption it follows that the induced cost of 1s that are to the left of R is at least 1 per every symbol 1 in Z_1LV_0 except from those 1s that are aligned with 1s in V_1 . The maximal number of 1s of Z_1LV_0 that can be aligned with 1s in V_1 is $|V_1| = l_2$. Thus, the total induced cost by 1s in Z_1LV_0 is $S_4 := \text{ones}(Z_1LV_0) - \text{ones}(V_1)$. Given that no symbol of D is aligned with a symbol of Z_1L , every symbol 0 in Z_1L induces a cost of 1. The total induced cost is $S_5 := \text{zeros}(Z_1L)$. We conclude that the total induced cost of symbols to the left of R in $VG_1(a)$ and to the left of D in $VG_2(b)$ is

$$\begin{aligned} S_6 &:= S_4 + S_5 = (\text{ones}(Z_1LV_0) - \text{ones}(V_1)) + \text{zeros}(Z_1L) \\ &= (\text{ones}(L) + |V_0| - |V_1|) + (|Z_1| + \text{zeros}(L)) = l + l_2. \end{aligned}$$

Case 2.3. A symbol 1 in V_2 is aligned with a symbol 1 in R . Suppose that this happens at position x . Then the induced cost of symbols 0 in R that are to the right of position x is at least $l_1 \geq dl_0 + d =: S_7$, which shows that the EDIT cost is at least

$$\begin{aligned} \text{EDIT}(VG(a), VG(b)) &\geq S_6 + S_1 + S_7 \\ &= (l + l_2) + l_2 + (dl_0 + d) = E_u. \end{aligned}$$

The last sub-case to consider complements cases 2.1, 2.2 and 2.3.

Case 2.4. The complement of cases 2.1, 2.2 and 2.3. This means that, if a symbol 1 in R (D , resp.) is aligned with a symbol 1, then the latter symbol comes from D (R , resp.).

First, observe that both R and D have d runs of 1s. Then we consider the following two sub-cases.

Case 2.4.1. There exists a run in R or D such that 1s in it are aligned with 1s in more than one run in the other sequence. In this case, the induced cost of symbols in R and D is at least $2l_1 \geq dl_0 + d =: S_8$, from which we conclude that

$$\begin{aligned} \text{EDIT}(\text{VG}_1(a), \text{VG}_2(b)) &\geq S_8 + S_6 + S_1 \\ &= (dl_0 + d) + (l + l_2) + l_2 = E_u. \end{aligned}$$

Case 2.4.2. The complement of case 2.4.1. Suppose that p runs in R have alignments with runs in D .

Case 2.4.2.1. $p = d$. In this case, for all $i = 1, 2, \dots, d$, the i -th run in R aligns with the i -th run in D . By the construction, the i -th run in R and D incur EDIT cost $\geq l_0$. Given that there exist at least one coordinate j such that $a_j \cdot b_j = 1$ (thus, the corresponding pair of runs incur cost $\geq 3l_0$), the total contribution of the runs of 1s in D and R to EDIT cost is at least $S_9 := l_0 \cdot (d - 1) + 3l_0$. Therefore,

$$\begin{aligned} \text{EDIT}(\text{VG}_1(a), \text{VG}_2(b)) &\geq S_6 + S_1 + S_9 \\ &= (l_2 + l) + l_2 + (l_0 \cdot (d - 1) + 3l_0) \geq 2l_2 + l + l_0d + d = E_u. \end{aligned}$$

Case 2.4.2.2. $p < d$. In this case, p pairs of runs of 1s incur total cost $\geq pl_0 =: S_{10}$. The rest of $d - p$ runs in R incur total cost $\geq (d - p)l_0 =: S_{11}$. The rest of $d - p$ runs in D incur total cost $\geq (d - p) \cdot 2l_0 =: S_{12}$.

$$\begin{aligned} \text{EDIT}(\text{VG}_1(a), \text{VG}_2(b)) &\geq S_6 + S_1 + S_{10} + S_{11} + S_{12} \\ &= (l_2 + l) + l_2 + pl_0 + (d - p)l_0 + (d - p) \cdot 2l_0 \geq 2l_2 + l + l_0d + d = E_u, \end{aligned}$$

where we use that $p < d$.

Case 3. Symbols of D are not aligned with symbols in RZ_2 . The analysis of the last main case is similar to the analysis of Case 2. We list the sub-cases below that we need to consider. Case 3.x or Case 3.x.y corresponds to Case 2.x or Case 2.x.y, respectively. We obtain the same lower bound on $\text{EDIT}(\text{VG}_1(a), \text{VG}_2(b))$ except that in this case the lower bound can be achieved.

Similarly as in Fact 2, we get that the total contribution of Z_1 and V_1 is at least $S'_1 = l_2$.

Case 3.1. A symbol 1 in V_0 is aligned with symbol 1 in V_1D . The last l_1 symbols 0 of D induce cost $l_1 \geq \text{ones}(R) + dl_0 + d =: S'_2$. The 0s to the right of V_0 induce cost $S'_3 = \text{zeros}(RZ_2)$. We get

$$\begin{aligned} \text{EDIT}(\text{VG}_1(a), \text{VG}_2(b)) &\geq S'_1 + S'_2 + S'_3 \\ &= l_2 + (\text{ones}(R) + dl_0 + d) + \text{zeros}(RZ_2) = E_u. \end{aligned}$$

Case 3.2. A symbol 1 in V_2 is aligned with symbol 1 in L . The last $l_1 \geq \text{ones}(R) + dl_0 + d =: S'_2$ zeros of L induce cost of at least S'_2 . Zeros to the right of V_0 induce cost at least $S'_3 = \text{zeros}(RZ_2)$. We get

$$\begin{aligned} \text{EDIT}(\text{VG}(a), \text{VG}(b)) &\geq S'_1 + S'_2 + S'_3 \\ &= l_2 + (\text{ones}(R) + dl_0 + d) + \text{zeros}(RZ_2) = E_u. \end{aligned}$$

We proceed by defining a quantity that we will need for the analysis of cases 3.3 and 3.4. The total induced cost by 1s in V_0RZ_2 is at least $S'_4 := \text{ones}(V_0RZ_2) - \text{ones}(V_2)$. The total induced cost

by 0s in RZ_2 is at least $S'_5 := \text{zeros}(RZ_2)$. The total induced cost by symbols to the right of L and to the right of D is at least

$$\begin{aligned} S'_6 &:= S'_4 + S'_5 = (\text{ones}(V_0 RZ_2) - \text{ones}(V_2)) + \text{zeros}(RZ_2) \\ &= (|V_0| + \text{ones}(R) - |V_2|) + (\text{zeros}(R) + |Z_2|) = l + l_2. \end{aligned}$$

Case 3.3. A symbol 1 in V_1 is aligned with a symbol 1 in L at position, say x . The induced cost of symbols 0 in L to the left of position x is at least $l_1/3 \geq dl_0 + d =: S'_7$. We get

$$\begin{aligned} \text{EDIT}(\text{VG}(a), \text{VG}(b)) &\geq S'_6 + S'_1 + S'_7 \\ &= (l + l_2) + l_2 + (dl_0 + d) = E_u. \end{aligned}$$

Case 3.4. The complement of cases 3.1, 3.2 and 3.3, i.e., if a symbol 1 in L (D) is aligned with symbol 1, then it comes from D (L).

We observe that L has $2d$ runs 1s (d runs of length 1 and d runs of length $3l_0$) and D has d runs of 1s.

Case 3.4.1. There exists a run in L or D such that 1s in it are aligned with 1s in more than one run in the other sequence. The induced cost by symbols in L and D is at least $l_1/3 \geq dl_0 + d =: S'_8$. We get

$$\begin{aligned} \text{EDIT}(\text{VG}_1(a), \text{VG}_2(b)) &\geq S'_8 + S'_6 + S'_1 \\ &= (dl_0 + d) + (l + l_2) + l_2 = E_u. \end{aligned}$$

Case 3.4.2. x runs of length 1 of ones and y runs of length $3l_0$ of ones from L align with runs in D . By construction, the x runs induce cost at least $S'_9 := x(2l_0 - 1)$, the y runs induce cost at least $S'_{10} := yl_0$. The remaining $d - x$ runs of length 1 induce cost at least $S'_{11} := (d - x) \cdot 1$, the remaining $d - y$ runs of length $3l_0$ induce cost at least $S'_{12} := (d - y) \cdot 3l_0$. We get

$$\begin{aligned} \text{EDIT}(\text{VG}_1(a), \text{VG}_2(b)) &\geq S'_6 + S'_1 + S'_9 + S'_{10} + S'_{11} + S'_{12} \\ &= (l_2 + l) + l_2 + x(2l_0 - 1) + yl_0 + (d - x) \cdot 1 + (d - y) \cdot 3l_0 \\ &\geq 2l_2 + l + l_0d + d = E_u, \end{aligned}$$

where the inequality can be made tight by setting $x = 0$ and $y = d$. □

3.2 The reduction

We proceed by concatenating vector gadgets into sequences.

We set $t = \max(|\text{VG}_1|, |\text{VG}_2|)$ and $T = 1000d \cdot t$. We define $\text{VG}'_k(a) = 2^T \text{VG}_k(a) 2^T$ for $k \in \{1, 2\}$. Let $f \in \{0, 1\}^d$ be a vector consisting of only 1s.

Let A and B be sets from the Orthogonal Vectors instance. We assume that $|A| \leq |B|$.

We define sequences

$$\begin{aligned} P_1 &= \bigcirc_{a \in A} \text{VG}'_1(a), \\ P_2 &= \left(\bigcirc_{i=1}^{|A|-1} \text{VG}'_2(f) \right) \left(\bigcirc_{b \in B} \text{VG}'_2(b) \right) \left(\bigcirc_{i=1}^{|A|-1} \text{VG}'_2(f) \right). \end{aligned}$$

Theorem 1. *Let $X := |A| \cdot E_u$. If there are two orthogonal vectors, one from set A , another from set B , then $\text{PATTERN}(P_1, P_2) \leq X - (E_u - E_s)$; otherwise we have $\text{PATTERN}(P_1, P_2) = X$.*

Proof. Follows from Lemmas 3 and 4 below. \square

Lemma 3. *If there are two orthogonal vectors, one from A , another from B , then*

$$\text{PATTERN}(P_1, P_2) \leq X - (E_u - E_s).$$

Proof. Let $a \in A$ and $b \in B$ be vectors such that $a \cdot b = 0$.

We can choose a contiguous subsequence s of P_2 consisting of a sequence of $|A_1|$ vector gadgets VG_2 such that s has the following property: transforming vector gadgets VG'_1 from P_1 and their corresponding vector gadgets VG'_2 from s into the same sequence one by one as per Lemma 2 and 1, we achieve cost smaller than the upper bound. We use the fact that at least one transformation is cheap because $a \cdot b = 0$ and we choose s so that $\text{VG}_1(a)$ and $\text{VG}_2(b)$ gets transformed into the same sequence. \square

Lemma 4. *If there are no two orthogonal vectors, one from A , another from B , then*

$$\text{PATTERN}(P_1, P_2) = X.$$

Proof. Consider a graph $(X_1 \cup X_2, E)$ with vertices $x_1(a) \in X_1, a \in A, x_2(b) \in X_2, b \in B$. Consider an optimal transformation of P_1 and a subsequence of P_2 into the same sequence according to Definition 1. We connect two vertices $x_1(a)$ and $x_2(b)$, iff $\text{VG}_1(a)$ and $\text{VG}_2(b)$ have an alignment in the transformation.

We want to claim that every vector gadget $\text{VG}_1(a)$ from P_1 contributes a cost of at least E_u to the final cost of $\text{PATTERN}(P_1, P_2)$. This will give $\text{PATTERN}(P_1, P_2) \geq X$. We consider connected components of the graph. We will show that a connected component that has $r \geq 1$ vertices from X_1 , contribute $\geq r \cdot E_u$ to the final cost of $\text{PATTERN}(P_1, P_2)$. From the case analysis below we will see that these contributions for different connected components are separate. Therefore, by summing up the contributions for all the connected components, we get $\text{PATTERN}(P_1, P_2) \geq |A| \cdot E_u = X$.

Consider a connected component of the graph with at least one vertex from X_1 . We examine several cases.

Case 1. The connected component has only one vertex from X_1 . Let $x_1(a)$ be the vertex.

Case 1.1. $x_1(a)$ is connected to more than one vertex. In this case, $\text{VG}_1(a)$ induces a cost of at least $2T > E_u$ (this cost is induced by symbols 2).

Case 1.2. $x_1(a)$ (corresponding to vector gadget $\text{VG}_1(a)$) is connected to only one vertex $x_2(b)$ (corresponding to vector gadget $\text{VG}_2(b)$). In this case, we can assume that no symbol of sequences $\text{VG}_1(a)$ and $\text{VG}_2(b)$ is aligned with symbol 2. If there is a symbol that gets aligned with symbol 2, we delete the corresponding symbol from $\text{VG}_1(a)$ or $\text{VG}_2(b)$. We notice that this replacement of substitution by deletion does not change the contribution of $\text{VG}_1(a)$ and $\text{VG}_2(b)$ to $\text{PATTERN}(P_1, P_2)$.

Let x be a contiguous substring of P_2 that achieves the minimum of $\text{EDIT}(P_1, x)$ (see Definition 1).

Case 1.2.1. The vector gadget $\text{VG}_2(b)$ is fully contained in the substring x . The lower bound on contribution of $\text{VG}_1(a)$ follows from Lemma 2.

Case 1.2.2. The complement of Case 1.2.1. We need to consider this case because of the following reason. We could potentially achieve a smaller than E_u contribution of $\text{VG}_1(a)$ to $\text{PATTERN}(P_1, P_2)$ by transforming $\text{VG}_1(a)$ and a *contiguous substring* of $\text{VG}_2(b)$ into the same string (instead of transforming $\text{VG}_1(a)$ and $\text{VG}_2(b)$ into the same string). In the next paragraph we show that this can't happen.

$\text{VG}_2(b)$ shares symbols with x and is not fully contained in x . $\text{VG}_2(b)$ must be the left-most (right-most, resp.) vector gadget in x but then T left-most (right-most, resp.) symbols 2 of $\text{VG}'_1(a)$ induce a cost of at least $T > E_u$.

Case 1.3. $x_1(a)$ is connected to no vertex. We get that $\text{VG}_1(a)$ induces cost of at least $|\text{VG}_1| > E_u$.

Case 2. The connected component has $r > 1$ vertices $x_1(a)$ from X_1 . In this case, the cost induced by the vector gadgets $\text{VG}_1(a)$ corresponding to the vertices from X_1 in the connected component is at least $(r - 1) \cdot 2T > r \cdot E_u$ (this cost is induced by symbols 2).

This finishes the argument that $\text{PATTERN}(P_1, P_2) \geq X$.

It remains to argue that we can achieve cost X (to show that $\text{PATTERN}(P_1, P_2) \leq X$) and it can be done as in Lemma 3. \square

A note on the vector gadget construction At the end of Section 3.1 we claim that $\text{EDIT}(\bar{\text{VG}}_1(a), \bar{\text{VG}}_2(b)) = t_0 + t \cdot (a \cdot b)$ and that this vector gadget construction does not and cannot possibly yield the desired reduction. The reason is as follows. If the edit distance between the vector gadgets grew with the inner product, then the large contribution from the “very” non-orthogonal vectors could overwhelm the small contribution from the pair of orthogonal vectors. As a result, we could not be able to tell whether two orthogonal vectors exist.

Furthermore, there is a more fundamental reason why the simple gadget does not work: the edit distance transformation corresponding to the cost $\text{EDIT}(\bar{\text{VG}}_1(a), \bar{\text{VG}}_2(b)) = t_0 + t \cdot (a \cdot b)$ uses only substitutions. The remainder of the hardness proof for PATTERN uses only substitutions as well. Therefore, if the construction worked, this would imply hardness of pattern matching with substitutions only. However, the latter problem can be solved in time $O(n \log n)$ by using Fast Fourier Transform, so the reduction cannot be valid.

3.3 Reducing PATTERN to EDIT

Let $|A| \leq |B|$. We set $P'_2 := P_2$ and $P'_1 := 3^{|P'_2|} P_1 3^{|P'_2|}$.

Theorem 2. *Let $Y := 2 \cdot |P'_2| + |A| \cdot E_u$. If there are no two orthogonal vectors, then $\text{EDIT}(P'_1, P'_2) = Y$; otherwise $\text{EDIT}(P'_1, P'_2) \leq Y - (E_u - E_s)$.*

Proof. Follows from Lemmas 5 and 6 below. \square

Lemma 5. *If there are two orthogonal vectors, then*

$$\text{EDIT}(P'_1, P'_2) \leq Y - (E_u - E_s).$$

Proof. We transform P_1 and a subsequence of P'_2 into the same sequence as in Lemma 3. We replace the remaining prefix and suffix of P'_2 with symbols 3 and delete the excess of symbols 3 from P'_1 . \square

Lemma 6. *If there are no two orthogonal vectors, then*

$$\text{EDIT}(P'_1, P'_2) = Y.$$

Proof. We can easily check that $\text{EDIT}(P'_1, P'_2) \leq Y$ in a way that is similar how it is done in Lemma 5. It remains to prove the opposite inequality.

P'_1 contains $2|P'_2|$ symbols 3. Those will incur a cost of at least $2|P'_2|$. P'_1 has the remaining subsequence P_1 , which will incur cost at least $\text{PATTERN}(P_1, P'_2)$. Using Lemma 4, we finish the proof. \square

As a result of all this, we get the following theorem.

Theorem 3. *If EDIT can be computed in time $O(n^{2-\delta})$ for some $\delta > 0$ on two sequences of length n over an alphabet of size 4, then Orthogonal Vectors Problem with $|A| = |B| = N$ and $A, B \subseteq \{0, 1\}^d$ can be solved in time $d^{O(1)} \cdot N^{2-\delta}$.*

Proof. The proof follows immediately from Theorem 2 by having $|A| = |B|$. \square

4 Acknowledgments

The authors thank Amir Abboud, Karl Bringmann, Sepideh Mahabadi and Ludwig Schmidt for providing helpful comments on an earlier version of the paper. This work was supported by grants from the NSF, the MADALGO center, and the Simons Investigator award.

References

- [ABW15] Amir Abboud, Arturs Backurs, and Virginia Vassilevska Williams. Quadratic-time hardness of lcs and other sequence similarity measures. *arXiv preprint arXiv:1501.07053*, 2015.
- [AKO10] Alexandr Andoni, Robert Krauthgamer, and Krzysztof Onak. Polylogarithmic approximation for edit distance and the asymmetric query complexity. In *FOCS*, pages 377–386. IEEE, 2010.
- [AW14] A. Abboud and V. Vassilevska Williams. Popular conjectures imply strong lower bounds for dynamic problems. *FOCS*, 2014.
- [AWW14] Amir Abboud, V. Vassilevska Williams, and Oren Weimann. Consequences of faster sequence alignment. *ICALP*, 2014.
- [BK15] Karl Bringmann and Marvin Künnemann. Quadratic conditional lower bounds for string problems and dynamic time warping. *arXiv preprint arXiv:1502.01063*, 2015.
- [Bri14] Karl Bringmann. Why walking the dog takes time: Frechet distance has no strongly subquadratic algorithms unless seth fails. *FOCS*, 2014.
- [IP01] Russell Impagliazzo and Ramamohan Paturi. On the complexity of k-sat. *Journal of Computer and System Sciences*, 62(2):367–375, 2001.
- [IPZ01] Russell Impagliazzo, Ramamohan Paturi, and Francis Zane. Which problems have strongly exponential complexity? *Journal of Computer and System Sciences*, 63:512–530, 2001.

- [MP80] William J Masek and Michael S Paterson. A faster algorithm computing string edit distances. *Journal of Computer and System sciences*, 20(1):18–31, 1980.
- [Nav01] Gonzalo Navarro. A guided tour to approximate string matching. *ACM computing surveys (CSUR)*, 33(1):31–88, 2001.
- [PW10] M. Patrascu and R. Williams. On the possibility of faster sat algorithms. *SODA*, 2010.
- [RW13] L. Roditty and V. Vassilevska Williams. Fast approximation algorithms for the diameter and radius of sparse graphs. *STOC*, 2013.
- [Wil05] Ryan Williams. A new algorithm for optimal 2-constraint satisfaction and its implications. *Theoretical Computer Science*, 348(2):357–365, 2005.