

Indexed Dynamic Programming to boost Edit Distance and LCSS Computation^{*}

Jérémy Barbay, Andrés Olivares

Departamento de Ciencias de la Computación, University of Chile, jeremy@barbay.cl, aolivare@dcc.uchile.cl

Abstract. There are efficient dynamic programming solutions to the computation of the Edit Distance from $S \in [1..\sigma]^n$ to $T \in [1..\sigma]^m$, for many natural subsets of edit operations, typically in time within $O(nm)$ in the worst-case over strings of respective lengths n and m (which is likely to be optimal), and in time within $O(n+m)$ in some special cases (e.g. disjoint alphabets). We describe how indexing the strings (in linear time), and using such an index to refine the recurrence formulas underlying the dynamic programs, yield faster algorithms in a variety of models, on a continuum of classes of instances of intermediate difficulty between the worst and the best case, thus refining the analysis beyond the worst case analysis. As a side result, we describe similar properties for the computation of the Longest Common Sub Sequence $\text{LCSS}(S, T)$ between S and T , since it is a particular case of Edit Distance, and we discuss the application of similar algorithmic and analysis techniques for other dynamic programming solutions. More formally, we propose a parameterized analysis of the computational complexity of the Edit Distance for various set of operators and of the Longest Common Sub Sequence in function of the area of the dynamic program matrix relevant to the computation.

1 Introduction

Given a set of edition operators on strings, a source string $S \in [1..\sigma]^n$ and a target string $T \in [1..\sigma]^m$ of respective lengths n and m on the alphabet $[1..\sigma]$, the EDIT DISTANCE is the minimum number of such operations required to transform the string S into the string T . Introduced in 1974 by Wagner and Fischer [16], such computation is a fundamental problem in Computer Science, with a wide range of applications, from text processing and information retrieval to computational biology. The typical edition distance between two strings is defined by the minimum number of **insertions**, **deletions** (in both cases, of a character at an arbitrary position of S) and **replacement** (of one character of S by some other) needed to transform the string S into T . Many generalizations have been defined in the literature, including weighted costs for the edit operations, and different sets of edit operations – the standard set is $\{\text{insertion}, \text{deletion}, \text{replacement}\}$.

Each distinct set of correction operators yields a distinct correction distance on strings (see Figure 1 for a summary). For instance, Wagner and Fischer [16] showed that for the three following operations, the **insertion** of a symbol at some arbitrary position, the **deletion** of a symbol at some arbitrary position, and the **replacement** of a symbol at some arbitrary position, the EDIT DISTANCE can be computed in time within $O(nm)$ and space within $O(n+m)$ using traditional dynamic programming techniques. As another variant of interest, Wagner and Lowrance [17] introduced the **Swap** operator (**S**), which exchanges the positions of two contiguous symbols. When considering only the **Swap** operator, one basically searches for the permutation transforming the source string S into the target string T : some adaptive sorting technique yields a minor improvement on the computation of the SWAP EDIT DISTANCE (see appendix A). For two of the newly defined distances, the INSERT SWAP EDIT DISTANCE and the DELETE SWAP EDIT DISTANCE (equivalent by symmetry), the best known algorithms take time exponential in the input size [4,5], which is likely to be optimal [16]. The EDIT DISTANCE itself is linked to many other problems: for instance, given the two same strings $S \in [1..\sigma]^n$ and $T \in [1..\sigma]^m$, the computation of the LONGEST COMMON SUB-SEQUENCE (LCSS) L between S and T is equivalent to the computation of the DELETE INSERT EDIT DISTANCE d , as the symbols deleted from S and inserted from T in order to “edit” S into T are exactly the same as the symbols deleted

^{*} Supported by project Fondecyt Regular no 1170366 from Conicyt.

Operators	(n, m) -Worst Case Complexity	Finer Results	
		Distance	Parikh vectors
Delete	$O(m)$ [15]	$O(n(1 + \lg d/n) \lg \lg \sigma)$ $O(n(1 + \lg d/n) \lg n)$	DNA
Insert	$O(n)$ [15]		
Replace	$O(n)$ [15]		
Swap	$O(n^2)$ [15]		
Delete, Insert	$O(nm)$ [7]	$O(d^2)$	$O(\sigma^2 nm \gamma^{\sigma-1})$ [4] $O(\sigma^2 nm \gamma^{\sigma-1})$ [4]
Delete, Replace	$O(nm)$ [18]	$O(d^2)$	
= Insert, Replace	$O(nm)$ [18]	$O(d^2)$	
Delete, Swap	NP-complete [16]	$O(1.6181^d m)$ [1]	
= Insert, Swap	NP-complete [16]	$O(1.6181^d m)$ [1]	
Replace, Swap	$O(nm)$ [18]	$O(d^2)$	
Delete, Insert, Replace	$O(nm)$ [7]	$O(d^2)$	
Delete, Insert, Swap	$O(nm)$ [16]	$O(d^2)$	
Delete, Replace, Swap	$O(nm)$ [16]	$O(d^2)$	
Insert, Replace, Swap	$O(nm)$ [16]	$O(d^2)$	
Delete, Insert, Replace, Swap	$O(nm)$ [7]	$O(d^2)$	

Fig. 1. Summary of results for various combinations of operators from the basic set $\{\text{Insert, Delete, Replace, Swap}\}$. The column labeled “ (n, m) -Worst Case Complexity” presents results in the worst case over instances of fixed sizes n and m , while the columns labeled “Finer Results” present results where the analysis was refined by various parameters: the distance d , the size σ of the alphabet, and some form of imbalance $\gamma = \max_{\alpha \in [1..\sigma]} \min\{n_\alpha, m_\alpha - n_\alpha\}$ between the **Parikh vectors** of S and T . For brevity, the only distance based on a single operator presented is the SWAP EDIT DISTANCE, as the computation of the others is always linear in the size of the input.

from S and T in order to produce L . Hence, the LCSS between S and T can be computed in time within $O(nm)$ and space within $O(n+m)$ using traditional dynamic programming techniques.

Most of these computational complexities are likely to be optimal in the worst case over instances of size (n, m) : the algorithms computing the three basic distances (INSERT EDIT DISTANCE, DELETE EDIT DISTANCE and REPLACE EDIT DISTANCE) in linear time are optimal as any algorithm must read the whole strings; the INSERT SWAP EDIT DISTANCE and its symmetric the DELETE SWAP EDIT DISTANCE are NP-hard to compute [18]; and in 2015 Backurs and Indyk [2] showed that the $O(n^2)$ upper bound for the computation of the DELETE INSERT REPLACE EDIT DISTANCE is optimal unless the *Strong Exponential Time Hypothesis* (SETH) is false.

More recently, Barbay and Pérez-Lantero [4,5], complementing Meister’s previous results [13] by the use of an index supporting the operators **rank** and **select** on strings, described an algorithm computing this distance in time within $O(\sigma^2 nm \gamma^{\sigma-1})$ in the worst case over instances where σ, n, m and γ are fixed, where $\gamma = \max_{\alpha \in [1..\sigma]} \min\{n_\alpha, m_\alpha - n_\alpha\}$ measures a form of imbalance between the frequency distributions of each string.

Hypothesis: Given this situation, is it possible to **take advantage of** indexing techniques supporting **rank** and **select** in order to **speed up the computation of** other edit distances? Can a similar analysis to that of Barbay and Pérez-Lantero [4,5] be applied to other edit distances? **Are there instances for which the edit distance is easier to compute, and do such instances occur in real applications** of the computation of the edit distance?

Our Results: We answer all those questions positively, and describe general techniques to refine the analysis of dynamic programs beyond the traditional analysis in the worst case over input of fixed size. More specifically, we analyze the computational cost of four EDIT DISTANCES using various **rank** and **select** text indices, in function of the **Parikh vector** [20] of the source S and target T strings. As a side result, this yields similar

properties for the computation of the LONGEST COMMON SUB SEQUENCE $\text{LCSS}(S, T)$ between S and T , as it can be deduced from the DELETE INSERT EDIT DISTANCE ($\text{LCSS}(S, T) = |S| + |T| - 2d_{DI}(S, T)$), and definitions and techniques which can be applied to other dynamic programs. After defining formally the notion of **Parikh's vector** and various index data structures supporting **rank** and **select** on strings in Section 2, we describe the algorithms taking advantage of such techniques in Section 3: for the LONGEST COMMON SUB SEQUENCE and DELETE-INSERT EDIT DISTANCE (Section 3.1), the DELETE INSERT REPLACE EDIT DISTANCE (Section 3.2), and finally for the DELETE-REPLACE EDIT DISTANCE and its dual the INSERT-REPLACE EDIT DISTANCE (Section 3.3). We describe some preliminary experiments and their results, which seem to indicate that those instances are not totally artificial and occur naturally in practical applications in Section 4. We conclude in Section 5 with a discussion of other potential refinement of the analysis, and the extension of our results to other EDIT DISTANCES.

2 Preliminaries

Before describing our proposed algorithms to compute various EDIT DISTANCES, we describe formally in Section 2.1 the notion of **Parikh vector** which is essential to our analysis technique; and in Section 2.2 two key implementations of indices supporting the **rank** and **select** operators on strings.

2.1 Parikh vector

Given positive integers σ and n , a string $S \in [1..\sigma]^n$, and the integers n_1, \dots, n_σ such that n_α denotes the number of occurrences of the letter $\alpha \in [1..\sigma]$ in the string S , the **Parikh vector** of S is defined [20] as $p(S) = (n_1, \dots, n_\sigma)$.

Barbay and Pérez-Lantero [4] refined the analysis of the INSERT SWAP EDIT DISTANCE from a string $S \in [1..\sigma]^n$ to a string $T \in [1..\sigma]^m$ via a function of the **Parikh vectors** (n_1, \dots, n_σ) of S and (m_1, \dots, m_σ) of T , the local imbalance $\gamma_\alpha = \min\{n_\alpha, m_\alpha - n_\alpha\}$ for each symbol $\alpha \in [1..\sigma]$, projected to a global measure of imbalance, $\gamma = \max_{\alpha \in [1..\sigma]} \gamma_\alpha$. In the worst case among instances of fixed **Parikh vector**, they describe an algorithm to compute the INSERT SWAP EDIT DISTANCE in time within

$$O\left(dn + d^2n \cdot \left(\sum_{\alpha=1}^{\sigma} (m_\alpha - \gamma_\alpha)\right) \cdot \prod_{\alpha \in \overline{[1..\sigma]}} (\gamma_\alpha + 1)\right),$$

where $\overline{[1..\sigma]} = \{\alpha \in [1..\sigma] : \gamma_\alpha > 0\}$ if $\prod_{\alpha \in [1..\sigma]} \gamma_\alpha = 0$, and $\overline{[1..\sigma]} = [1..\sigma] \setminus \{\arg \min_{\alpha \in [1..\sigma]} \gamma_\alpha\}$ otherwise. This formula simplifies to within $O(\sigma^2 nm \gamma^{\sigma-1})$ in the worst case over instances where σ, n, m and γ are fixed.

Such a vector is essential to the fine analysis of dynamic programs for computing EDIT DISTANCES when using operators whose running time depends on the number of occurrence of each symbol, such as for the **rank** and **select** operators described in the next section.

2.2 Rank and Select in Strings

For every string $X \in \{S, L\}$ and integer $i \in [1..|X|]$, $X[i]$ denotes the i -th symbol of X from left to right. For every pair of integers $i, j \in [1..|X|]$ such that $i \leq j$, $X[i..j]$ denotes the substring of X from the i -th symbol to the j -th symbol, and for every pair of integers $i, j \in [1..|X|]$ such that $j < i$, $X[i..j]$ denotes the empty string.

Given a symbol $\alpha \in [1..\sigma]$, an integer $i \in [1..|X|]$ and an integer $k > 0$, the operation $\text{rank}(X, i, \alpha)$ denotes the number of occurrences of the symbol α in the substring $X[1..i]$, and the operation $\text{select}(X, k, \alpha)$ denotes the value $j \in [1..|X|]$ such that the k -th occurrence of α in X is precisely at position j , if j exists. If j does not exist, then $\text{select}(X, k, \alpha)$ is *null*.

A simple way to support the **rank** and **select** operators in reasonably good time consists in, for each symbol $\alpha \in [1..\sigma]$, listing all the occurrences of α in a sorted array (called a “Posting List” [21]): supporting the **select** operator reduces to a simple access to the sorted array corresponding to the symbol α ; while supporting the **rank** operator reduces to a SORTED SEARCH in the same array, which can be simply implemented by a **Binary Search**, or more efficiently in practice by a **Doubling Search** [6] in time within $O(q_\alpha \lg(n_\alpha/q_\alpha))$ when supporting q_α monotone queries in a posting list of size n_α (for a given symbol $\alpha \in [1..\sigma]$).

Lemma 1. *Given a string $S \in [1..\sigma]^n$ of **Parikh vector** (n_1, \dots, n_σ) , there exists an index using $n + \sigma$ machine words, which can be computed in time linear in the size n of S in order to support the operators **rank** and **select** in time within $O(q_\alpha \lg(n_\alpha/q_\alpha))$ in the comparison based decision tree model, when q_α of those queries concern the symbol $\alpha \in [1..\sigma]$.*

Golynski *et al.* [10] described a more sophisticated (but asymptotically more efficient) way to support the **rank** and **select** operators in the RAM model, via a clever reduction to Y-Fast Trees on permutations supporting the operators in time within $O(\lg \lg \sigma)$. Barbay *et al.* [3] showed that it can be done on a compressed representation of the text.

Lemma 2. *Given a string $S \in [1..\sigma]^n$, there exists an index using space within $o(n \lg \sigma)$ bits, which can be computed in time linear in the size n of S in order to support the operators **rank** and **select** in time within $O(\lg \lg \sigma)$ in the RAM model.*

We describe how to use those techniques to speed up the computation of various EDIT DISTANCES in the following sections.

3 Adaptive Dynamic Programs

For each of the problems considered, we describe how to compute a subset of the values computed by classical dynamic programs. We start with the computation of the LONGEST COMMON SUB SEQUENCE (LCSS) and the DELETE INSERT (DI) EDIT DISTANCE (Section 3.1) because it is the simplest; extend its results to the computation of the LEVENSHTAIN EDIT DISTANCE (Section 3.2); and project those to the computation of the DELETE REPLACE (DR) EDIT DISTANCE and its symmetric INSERT REPLACE (IR) EDIT DISTANCE (Section 3.3).

3.1 LCSS and DI-Edit Distance

The DELETE INSERT EDIT DISTANCE is a classical problem in Stringology [7], if only as a variant of the LONGEST COMMON SUB SEQUENCE problem. It is classically computed using dynamic programming: we describe the classical solution first, which we then refine in a simplistic way, as a pedagogical introduction to a more sophisticated refinement.

Classical solution: Given two strings $S \in [1..\sigma]^n$ and $T \in [1..\sigma]^m$, we note $d_{DI}(n, m)$ the DELETE INSERT EDIT DISTANCE from S to T . If the last symbols of S and T match, the edit distance is the same as the edit distance between the prefixes of respective lengths $n - 1$ and $m - 1$ of S and T . Otherwise, the edit distance is the minimum between the edit distance when inserting a copy of the last symbol of T in S (i.e. deleting this symbol in T) and the edit distance when deleting the mismatching symbol in T . More formally:

$$d_{DI}(S[1..n], T[1..m]) = \begin{cases} n & \text{if } m = 0; \\ m & \text{if } n = 0; \\ d_{DI}(S[1..n-1], T[1..m-1]) & \text{if } S[n] = T[m]; \text{ and} \\ 1 + \min \left\{ \begin{array}{l} d_{DI}(S[1..n-1], T[1..m]), \\ d_{DI}(S[1..n], T[1..m-1]) \end{array} \right\} & \text{otherwise.} \end{cases}$$

This recursive definition directly yields an algorithm to compute the DELETE INSERT EDIT DISTANCE from S to T in time within $O(nm)$ and space within $O(n + m)$. We describe in the next section a technique taking advantage of the discrepancies between the **Parikh vectors** of S and T .

A Pedagogical Example: Given two strings $S \in [1..\sigma]^n$ and $T \in [1..\sigma]^m$, for each symbol $\alpha \in [1..\sigma]$, let's note n_α and m_α the number of occurrences of α respectively in S and T . Assembled in a vector, those form the **Parikh vectors** $(n_\alpha)_{\alpha \in [1..\sigma]}$ for S and $(m_\alpha)_{\alpha \in [1..\sigma]}$ for T . Barbay and Pérez-Lantero [4] described an algorithm to compute the INSERT SWAP EDIT DISTANCE which complexity is expressed in function of how the **Parikh vectors** of S and T differ. Likewise, we describe how those affect the difficulty of computing the DELETE INSERT EDIT DISTANCE from S to T .

Consider in Figure 2 the graphical representation $M_{DI}(S, T)$ of the dynamic program computing the DELETE INSERT EDIT DISTANCE from S to T , following the dynamic program described in the previous section. For general $i \in [1..n]$ and $j \in [1..m]$, the i -th value in the j -th row, $a = M_{DI}(S, T)[i, j] = d_{DI}(S[1..i], T[1..j])$ is computed by taking the minimum between $b = M_{DI}(S, T)[i-1, j] = d_{DI}(S[1..i-1], T[1..j])$ and $c = M_{DI}(S, T)[i, j-1] = d_{DI}(S[1..i], T[1..j-1])$, the value directly on the left and directly above it: $a = \min\{b, c\}$.

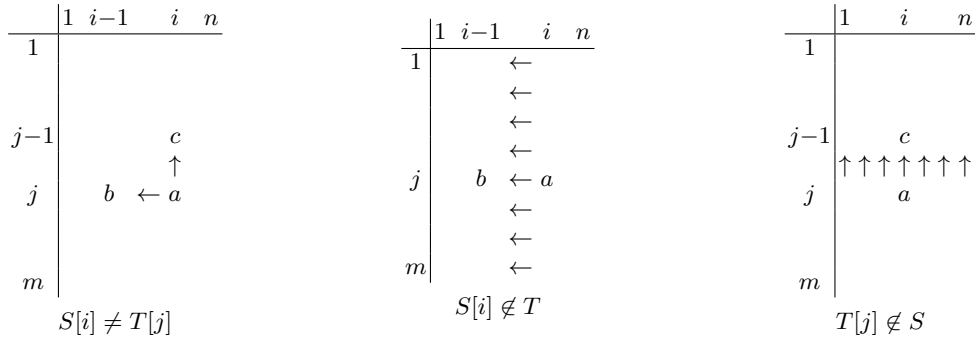


Fig. 2. A graphical representation of the dynamic program computing the DELETE INSERT EDIT DISTANCE from S to T in the general case ($S[i] \neq T[j]$) and in the particular case where the symbol at position i in S does not occur in T ($S[i] \notin T$), and where the symbol at position j in T does not occur in S ($T[j] \notin S$).

Consider a particular position $i \in [1..n]$ in S such that the symbol $S[i]$ at this position does not occur in T (i.e. $S[i] \notin T$): this symbol will be deleted in any edition of S into T , so that each value in the column i can be obtained by merely duplicating the corresponding one in the column $i-1$. Similarly, consider a particular position $j \in [1..m]$ in T such that the symbol $T[j]$ at this position does not occur in S (i.e. $T[j] \notin S$): this symbol will be inserted in any edition of S into T , so that each value in the row j can be obtained by merely duplicating the corresponding one in the row $j-1$. The duplication of such columns and row can be simulated in constant time during the execution of the dynamic program, thus reducing the complexity to within $O(n'm' + n + m + \sigma)$ where n' and m' are the lengths of S and T once projected to the intersection of their effective alphabets: $n' = \sum_{\alpha, m_\alpha > 0} n_\alpha$ and $m' = \sum_{\alpha, n_\alpha > 0} m_\alpha$. We show in the next section how to further refine this technique, in order to take advantage of rare symbols in each string.

Refined Analysis: We described in the previous section how to take advantage of the fact that some elements appear in one string, but not in the other. It is natural to wonder if a similar technique can take advantage of cases where a symbol occurs many time in one string, but occurs only once in the other: at some point, the dynamic program will reduce to the case described in the previous section. To be able to notice when this happens, one would need to maintain dynamically the counters of occurrences of each symbol during the execution of the dynamic program, or more simply pre-compute an index on S and T supporting the operators **rank** and **select** on it.

Given the support for the **rank** and **select** operators on both S and T , we can refine the dynamic program to compute the distance $d_{DI}(n, m)$ as follows:

$$d_{DI}(n, m) = \begin{cases} n & \text{if } m == 0; \\ m & \text{if } n == 0; \\ d_{DI}(n-1, m-1) & \text{if } S[n] == T[m]; \\ 1 + d_{DI}(n-1, m) & \text{if } \mathbf{rank}(S, T[m]) == 0; \\ 1 + d_{DI}(n, m-1) & \text{if } \mathbf{rank}(T, S[n]) == 0; \\ \min \left\{ \begin{array}{l} 1 + d_{DI}(n-1, m-1), \\ n - \mathbf{select}(S, T[m]) \\ \quad + d_{DI}(\mathbf{select}(S, T[m], \mathbf{rank}(S, T[m]) - 1) - 1, m-1), \\ m - \mathbf{select}(T, S[n]) \\ \quad + d_{DI}(n-1, \mathbf{select}(T, S[n], \mathbf{rank}(T, S[n]) - 1)) - 1 \end{array} \right\} & \text{otherwise.} \end{cases}$$

The running time of the algorithm can then be expressed in function of the number of recursive calls, the number of **rank** and **select** operations performed on the strings, in order to yield various running times depending upon the solution used to support the **rank** and **select** operators.

Theorem 1. *Given two strings $S \in [1..\sigma]^n$ and $T \in [1..\sigma]^m$ of respective **Parikh vectors** $(n_a)_{a \in [1..\sigma]}$ and $(m_a)_{a \in [1..\sigma]}$, the dynamic program above computes the DELETE INSERT EDIT DISTANCE from S to T and the LONGEST COMMON SUB SEQUENCE between S and T*

1. through at most $4 \sum_{a \in [1..\sigma]} n_a m_a$ recursive calls;
2. within $O(\sum_{a \in [1..\sigma]} n_a m_a)$ operations **rank** or **select**;
3. in time within $O(\sum_{a \in [1..\sigma]} n_a m_a \times \lg(\max_a \{n_a, m_a\}) \times \lg(nm))$ in the comparison model; and
4. in time within $O(\sum_{a \in [1..\sigma]} n_a m_a \times \lg \lg \sigma \times \lg(nm))$ in the RAM memory model;

Proof. We prove point (1) by an amortization argument. Point (2) is a direct consequence of point (1), given that each recursive call performs a finite number of calls to the **rank** and **select** operators. Point (3) is a simple combination of Point (2) with the classical *inverted posting list* implementation [21] of an index supporting the **select** operator in constant time and the **rank** operator via *doubling search* [6]; while point (4) is a simple combination of Point (2) with the index described by Golynski *et al.* [10] to support the **rank** and **select** operators.

Albeit quite simple, this results corresponds to real improvement in practice: see in Figure 3 how the number of recursive calls is reduced by using such indexes. Moreover, such a refinement of the analysis and optimization of the computation can be applied to more than the DELETE INSERT EDIT DISTANCE: in the next sections, we describe a similar one for computing the LEVENSHTAIN DISTANCE (Section 3.2) and the DELETE REPLACE and INSERT REPLACE EDIT DISTANCE (Section 3.3).

3.2 Levenshtein Distance, or DIR-Edit Distance

In information theory, linguistics and computer science, the LEVENSHTAIN DISTANCE is a string metric for measuring the difference between two sequences [7]. It generalizes the DELETE INSERT EDIT DISTANCE explored in the previous section by adding the **Replace** operator to the operators **Delete** and **Insert** (so that it can be also called the DELETE INSERT REPLACE EDIT DISTANCE, or *DIR* for short). The recursion traditionally used is a mere extension from the one described in the previous section:

$$d_{DIR}(n, m) = \begin{cases} m & \text{if } n == 0; \\ +\infty & \text{if } n > m; \\ d_{DIR}(n-1, m-1) & \text{if } S[n] == T[m]; \text{ and} \\ 1 + \min \left\{ \begin{array}{l} d_{DIR}(n, m-1), \\ d_{DIR}(n-1, m-1) \end{array} \right\} & \text{otherwise.} \end{cases}$$

The adaptive version is only a technical extension of the one for the DELETE INSERT EDIT DISTANCE:

$$\left\{ \begin{array}{l} n \text{ if } m == 0; \\ m \text{ if } n == 0; \\ d_{DIR}(n-1, m-1) \text{ if } S[n] == T[m]; \\ 1 + d_{DIR}(n-1, m-1) \text{ if } \mathbf{rank}(S, T[m]) == 0 \\ \text{and } \mathbf{rank}(T, S[n]) == 0 \text{ (REPLACE);} \\ 1 + d_{DIR}(n-1, m) \text{ if } \mathbf{rank}(S, T[m]) == 0 \\ \text{but } \mathbf{rank}(T, S[n]) > 0 \text{ (DELETE);} \\ 1 + d_{DIR}(n, m-1) \text{ if } \mathbf{rank}(T, S[n]) == 0 \\ \text{but } \mathbf{rank}(S, T[m]) > 0 \text{ (INSERT);} \\ \min \left\{ \begin{array}{l} n - \mathbf{select}(S, T[m]) \\ + d_{DIR}(\mathbf{select}(S, T[m], \mathbf{rank}(S, T[m]) - 1) - 1, m-1) \text{ (DELETE) ,} \\ m - \mathbf{select}(T, S[n]) \\ + d_{DIR}(n-1, \mathbf{select}(T, S[n], \mathbf{rank}(T, S[n]) - 1)) - 1 \text{ (INSERT) ,} \\ 1 + d_{DIR}(n-1, m-1) \text{ (REPLACE)} \end{array} \right\} \end{array} \right\} \text{ otherwise.}$$

The refined analysis yields similar results (we omit the proof for lack of space):

Theorem 2. *Given two strings $S \in [1..\sigma]^n$ and $T \in [1..\sigma]^m$ of respective **Parikh vectors** $(n_a)_{a \in [1..\sigma]}$ and $(m_a)_{a \in [1..\sigma]}$, the dynamic program above computes the LEVENSHTAIN EDIT DISTANCE from S to T*

1. through at most $4 \sum_{a \in [1..\sigma]} n_a m_a$ recursive calls;
2. within $O(\sum_{a \in [1..\sigma]} n_a m_a)$ operations **rank** or **select**;
3. in time within $O(\sum_{a \in [1..\sigma]} n_a m_a \times \lg(\max_a \{n_a, m_a\}) \times \lg(nm))$ in the comparison model; and
4. in time within $O(\sum_{a \in [1..\sigma]} n_a m_a \times \lg \lg \sigma \times \lg(nm))$ in the RAM memory model;

It is important to note that for two strings S and T , the computation of the LEVENSHTAIN EDIT DISTANCE from S to T actually generates more recursive calls than the computation of the DELETE INSERT EDIT DISTANCE from S to T , but that the analysis above does not capture this difference. In the following section, we project this result to two equivalent edit distances, the DELETE REPLACE and INSERT REPLACE EDIT DISTANCES, for which the dynamic program explores only half of the position in the dynamic program matrix compared to the LEVENSHTAIN EDIT DISTANCE or DELETE INSERT EDIT DISTANCE.

3.3 DR-Edit Distance and IR-Edit Distance

Given a source string $S \in [1..\sigma]^n$ and a target string $T \in [1..\sigma]^m$, the DELETE REPLACE EDIT DISTANCE from S to T and the INSERT-REPLACE EDIT DISTANCE from T to S are the same, as the sequence of **Insert** or **Replace** operations transforming S into T is the symmetric to the sequence of **Delete** or **Replace** operations transforming T back into S .

As before, if the last symbols of S and T match, the edit distance is the same as the edit distance between the prefixes of respective lengths $n-1$ and $m-1$ of S and T . Otherwise, the edit distance is the minimum between the edit distance when inserting a copy of the last symbol of T in S (i.e. deleting this symbol in T) and the edit distance when replacing the mismatching symbol in S by the corresponding one in T . More formally:

$$d_{DR}(S[1..n], T[1..m]) = \begin{cases} m \text{ if } n == 0; \\ +\infty \text{ if } n > m; \\ d_{DR}(S[1..n-1], T[1..m-1]) \text{ if } S[n] == T[m]; \text{ and} \\ 1 + \min \left\{ \begin{array}{l} d_{DR}(S[1..n], T[1..m-1]), \\ d_{DR}(S[1..n-1], T[1..m-1]) \end{array} \right\} \text{ otherwise.} \end{cases}$$

Using a few more optimizations than in Section 3.1, this recursive definition yields an algorithm to compute the INSERT REPLACE EDIT DISTANCE from S to T in time within $O(m^2)$ and space within $O(m)$.

One can observe a few optimizations, such as that the edit distance can be computed in time linear in m as soon as n is equal to m , as no further **Insert** operation can be performed.

As in the two previous sections, given the support for the **rank** and **select** operators on both S and T , we can refine the dynamic program to compute the DELETE REPLACE EDIT DISTANCE $d_{DR}(n, m)$ as follows:

$$\left\{ \begin{array}{l} n \text{ if } m == 0; \\ \infty \text{ if } n < m; \\ d_{DR}(n-1, m-1) \text{ if } S[n] == T[m]; \\ 1 + d_{DR}(n-1, m) \text{ if } \mathbf{rank}(T, S[n]) == 0 \text{ (DELETE)}; \\ 1 + d_{DR}(n-1, m-1) \text{ if } \mathbf{rank}(S, T[n]) == 0 \text{ (REPLACE)}; \\ \min \left\{ \begin{array}{l} n - \mathbf{select}(S, T[m]) \\ + d_{DR}(\mathbf{select}(S, T[m], \mathbf{rank}(S, T[m]) - 1), m-1) \text{ (DELETE)}, \\ 1 + d_{DR}(n-1, m-1) \text{ (REPLACE)} \end{array} \right\} \end{array} \right\} \text{ otherwise.}$$

The analysis from the two previous sections projects to a similar result.

Theorem 3. *Given two strings $S \in [1..\sigma]^n$ and $T \in [1..\sigma]^m$ of respective **Parikh vectors** $(n_a)_{a \in [1..\sigma]}$ and $(m_a)_{a \in [1..\sigma]}$, the dynamic program above computes the DELETE REPLACE EDIT DISTANCE from S to T*

1. *through at most $4 \sum_{a \in [1..\sigma]} n_a m_a$ recursive calls;*
2. *within $O(\sum_{a \in [1..\sigma]} n_a m_a)$ operations **rank** or **select**;*
3. *in time within $O(\sum_{a \in [1..\sigma]} n_a m_a \times \lg(\max_a \{n_a, m_a\}) \times \lg(nm))$ in the comparison model; and*
4. *in time within $O(\sum_{a \in [1..\sigma]} n_a m_a \times \lg \lg \sigma \times \lg(nm))$ in the RAM memory model;*

Parameterizing the analysis of the computation of the LONGEST COMMON SUB SEQUENCE, of the LEVENSHTAIN EDIT DISTANCE and of the DELETE REPLACE or INSERT REPLACE EDIT DISTANCE would be only of moderate theoretical interest, if it did not correspond to some correspondingly “easy” instances in practice. In the next section we describe some preliminary experimental results which seem to indicate the existence of such “easy” instances in information retrieval.

4 Experiments

In order to test the practicality of the parameterization and algorithms described in the previous section, we performed some preliminary experiments on some public data sets from the GUTENBERG project [12]. We describe the data set and experimental setup in Section 4.1, and the preliminary results and their interpretation in Section 4.2.

4.1 Data Sets

Started by Michael Hart in 1971 [19], the GUTENBERG project gathers electronic copies of public domain books, and as such is a publically available data set for testing algorithms on real text. We considered each text as a sequence of words (hence considering as equivalent all the word separations, from blank spaces to punctuations and line jumps), which results in large alphabets. Due to some problems with the implementation, we could not run the algorithms for texts larger than 32kB (a memory issue with a library in Python), so we extracted the first 32kB of the texts “Romeo & Juliet” (English), “Romeo & Julia” (German), “Hamlet” (German), and “Punch or the London chivalry vol 99” (English); the last text being a randomly picked non Shakespeare text.

Delete Insert Edit Distance Recursions

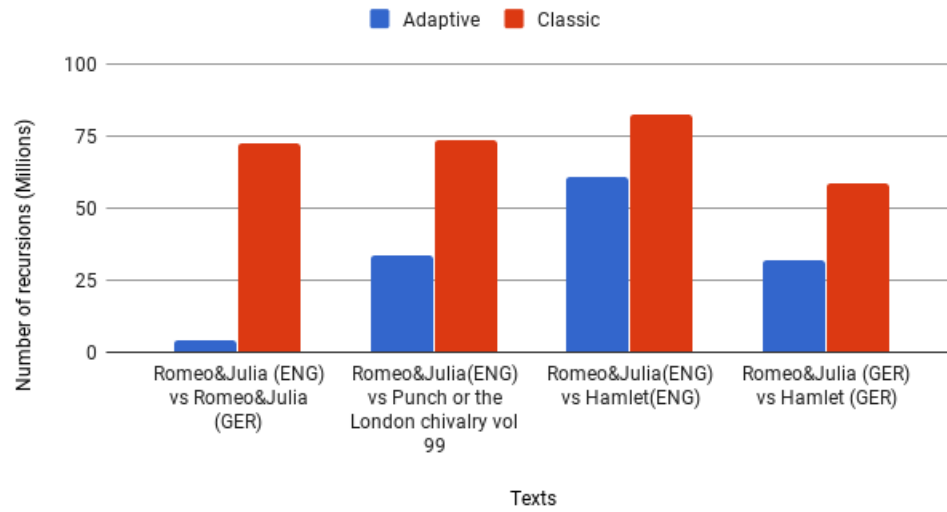


Fig. 3. Experimental results for the computation of the DELETE INSERT EDIT DISTANCE (and, by extension, of the computation of the LONGEST COMMON SUB SEQUENCE) by the adaptive algorithm described in Section 3.1.

Delete Insert Replace Edit Distance Recursions

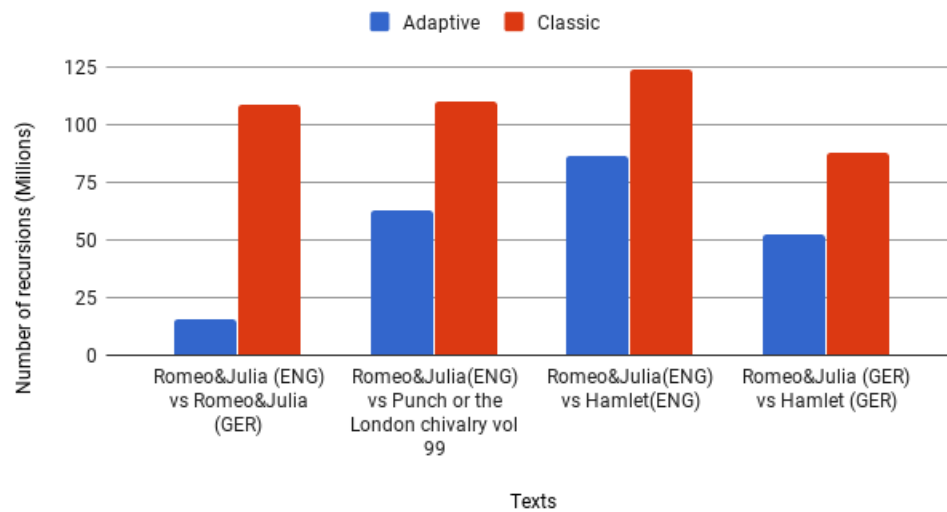
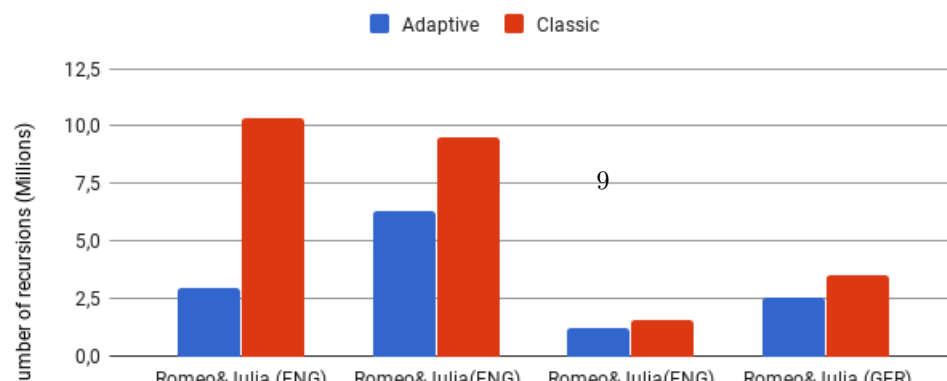


Fig. 4. Experimental results for the computation of the LEVENSHTAIN EDIT DISTANCE by the adaptive algorithm described in Section 3.2.

Delete Replace Edit Distance Recursions



4.2 Experimental Results

Figures 3, 4 and 5 show the number of recursive calls from the main recursive function for four pairs of texts for each algorithm described in Section 3: “Romeo & Juliet” (English) vs “Punch or the London chivalry vol 99” (English), “Romeo & Juliet” (English) vs “Romeo & Julia” (German), “Romeo & Juliet” (English) vs “Hamlet” (English), and “Romeo & Julia” (German) vs “Hamlet” (German).

For the three types of EDIT DISTANCES and the four pairs of texts, the adaptive variants perform less recursive calls. For the three types of EDIT DISTANCES, the difference in the number of recursive calls is less between the two texts from the same author (i.e. “Romeo & Juliet” (English) vs “Hamlet” (English) and “Romeo & Julia” (German) vs “Hamlet” (German)), because the vocabulary is the same, and is the most between texts of distinct languages (i.e. “Romeo & Juliet” (English) vs “Romeo & Julia” (German)), because the vocabulary (i.e. the alphabet) is mostly distinct. Still, for two texts in the same language, but from distinct authors (i.e. “Romeo & Juliet” (English) vs “Punch or the London chivalry vol 99” (English)), the difference is quite sensible.

Obviously, those experimental results are only preliminary, and a more thorough study is needed (and underway), both with a larger data set and with a larger range of measures, from the *running time* with various indexing data structures supporting the operators `rank` and `select`, to the number of entries of the dynamic program matrix being effectively computed. We discuss additional perspectives for future work in the next section.

5 Discussion

We have shown how the computation of other EDIT DISTANCES than the INSERT SWAP and DELETE SWAP EDIT DISTANCE is also sensitive to the **Parikh vectors** of the input. We discuss here various directions in which these results can be extended, from the possibility of proving conditional lower bounds in the refined analysis model, to further refinements of the analysis for these same EDIT DISTANCES, and to the analysis of other dynamic programs.

Adaptive Conditional Lower Bounds: Backurs and Indyk [2] showed that the $O(n^2)$ upper bound for the computation of the DELETE INSERT REPLACE EDIT DISTANCE is optimal unless the *Strong Exponential Time Hypothesis* (SETH) is false, and since then the technique has been applied to various other related problems. Should the reduction from the SETH to the EDIT DISTANCE computation be refined as shown here for the upper bound, it would speak in favor of the optimality of the analysis.

Other measures of difficulty: Abu-Khzam et al. [1] described an algorithm computing the INSERT SWAP EDIT DISTANCE d from S to T in time within $O(1.6181^d m)$, which is polynomial in the size of the input if exponential in the output size d . The output distance d itself can be as large as n , but such instances are not necessarily difficult: Barbay and Pérez-Lantero [4] showed that the gap vector between the **Parikh vectors** separate the hard instances from the easy one, and we showed that the same applied to other EDIT DISTANCES.

But still, among instances of fixed input size, output distance, and imbalance between the **Parikh vectors**, there are instances easier than others (e.g. the computation of the INSERT SWAP EDIT DISTANCE on an instance where all the **insertions** are in the left part of S while all the **swaps** are in the right part of S). A measure which to refine the analysis would be the cost of encoding a *certificate* of the EDIT DISTANCE, one which is easier to check than recomputing the distance itself.

Indexed Dynamic Programming: Our results are close in spirit to those in fixed-parameter complexity, but with an important difference, namely, trying to spot one or more parameters that explain what makes an instance hard or easy. For the computation of the INSERT SWAP and DELETE SWAP EDIT DISTANCES, the size of the alphabet d makes the difference between polynomial time and NP-hardness. However, Barbay and

Pérez-Lantero [4] showed that different instances of the same size can exhibit radically different costs-for a given fixed algorithm. The parameterized analysis captures in parameters the cause for such cost differences. We described how the same logic applies to other types of EDIT DISTANCES, and it is likely that similar situations happen with many other algorithms based on dynamic programming, such as the computation of the FRÉCHET DISTANCE [11], the DISCRETE FRÉCHET DISTANCE [9] and the decision problem ORTHOGONAL VECTOR [8].

Acknowledgments: The author would like to thank Pablo Pérez-Lantero for introducing the problem of computing the EDIT DISTANCE between strings; Felipe Lizama for a semester of very interesting discussions about this approach; and an anonymous referee from the journal Transaction on Algorithms for his positive feedback and encouragement. **Funding:** Jérémy Barbay is partially funded by the project Fondecyt Regular no. 1170366 from Conicyt. **Data and Material Availability:** The source of this article, along with the code and data used for the experiments described within, will be made publicly available upon publication at the url <https://github.com/FineGrainedAnalysis/EditDistances>.

References

1. Abu-Khzam, F.N., Fernau, H., Langston, M.A., Lee-Cultura, S., Stege, U.: Charge and reduce: A fixed-parameter algorithm for string-to-string correction. *Discrete Optimization (DO)* **8**(1), 41 – 49 (2011)
2. Backurs, A., Indyk, P.: Edit distance cannot be computed in strongly subquadratic time (unless SETH is false). In: *Proceedings of the annual ACM Symposium on Theory Of Computing (STOC)* (2015)
3. Barbay, J., Claude, F., Gagie, T., Navarro, G., Nekrich, Y.: Efficient fully-compressed sequence representations. *Algorithmica (ALGO)* **69**(1), 232–268 (2014)
4. Barbay, J., Pérez-Lantero, P.: Adaptive computation of the swap-insert correction distance. In: *Proceedings of the Annual Symposium on String Processing and Information Retrieval (SPIRE)*. pp. 21–32 (2015)
5. Barbay, J., Pérez-Lantero, P.: Adaptive computation of the swap-insert correction distance. In: *ACM Transactions on Algorithms (TALG)* (2018), accepted on [2018-05-25 Fri], to appear.
6. Bentley, J.L., Yao, A.C.C.: An almost optimal algorithm for unbounded searching. *Information Processing Letters (IPL)* **5**(3), 82–87 (1976)
7. Bergroth, L., Hakonen, H., Raita, T.: A survey of longest common subsequence algorithms. In: *Proceedings of the 11th Symposium on String Processing and Information Retrieval (SPIRE)*. pp. 39–48 (2000)
8. Bringmann, K.: Why walking the dog takes time: Fréchet distance has no strongly subquadratic algorithms unless SETH fails. In: *Proceedings of the 2014 IEEE 55th Annual Symposium on Foundations of Computer Science*. pp. 661–670. FOCS '14, IEEE Computer Society, Washington, DC, USA (2014)
9. Eiter, T., Mannila, H.: Computing discrete Fréchet distance. Tech. rep., Christian Doppler Labor für Expertensysteme, Technische Universität Wien (1994)
10. Golynski, A., Munro, J.I., Rao, S.S.: Rank/select operations on large alphabets: A tool for text indexing. In: *Proceedings of the Seventeenth Annual ACM-SIAM Symposium on Discrete Algorithm (SODA)*. pp. 368–373. SODA '06, Society for Industrial and Applied Mathematics, Philadelphia, PA, USA (2006)
11. H., A., M., G.: Computing the Fréchet distance between two polygonal curves. *International Journal of Computational Geometry and Applications (IJCGA)* **5**(1–2), 75–91 (1995)
12. Hart, M.: Gutenberg project. Online at <https://www.gutenberg.org/> (last accessed on [2018-05-27 Sun])
13. Meister, D.: Using swaps and deletes to make strings match. *Theoretical Computer Science (TCS)* **562**(0), 606 – 620 (2015)
14. Moffat, A., Petersson, O.: An overview of adaptive sorting. *Australian Computer Journal (ACJ)* **24**(2), 70–77 (1992)
15. Spreen, T.D.: The Binary String-to-String Correction Problem. Master’s thesis, University of Victoria, Canada (2013)
16. Wagner, R.A., Fischer, M.J.: The string-to-string correction problem. *Journal of the ACM (JACM)* **21**(1), 168–173 (1974)
17. Wagner, R.A., Lowrance, R.: An extension of the string-to-string correction problem. *Journal of the ACM (JACM)* **22**(2), 177–183 (1975)
18. Wagner, R.A.: On the complexity of the extended string-to-string correction problem. In: *Proceedings of the annual ACM Symposium on Theory Of Computing (STOC)*. pp. 218–223. STOC '75, ACM (1975)

19. Wikipedia: Project gutenber. Online at https://en.wikipedia.org/wiki/Project_Gutenberg (last accessed on [2018-05-27 Sun])
20. Wikipedia, **Website.**: Parikh's theorem, last accessed on 2017-05-08.
21. Witten, I.H., Moffat, A., Bell, T.C.: Managing gigabytes : compressing and indexing documents and images. The Morgan Kaufmann series in multimedia information and systems, San Francisco, Calif. Morgan Kaufmann Publishers (1999)

APPENDIX

In this appendix, we briefly discuss some minor topics, such as how the algorithm **Local Insertion Sort** described and analyzed by Moffat and Petersson [14] combined with an index supporting the **rank** and **select** operators potentially yields a faster computation of the SWAP EDIT DISTANCE (Section A), or how to combine techniques that take advantage of the **Parikh vectors** of the input strings with techniques that take advantage of the output distance (Section B).

A Adaptive Computation of the SWAP EDIT DISTANCE

Out of the $2^4 - 1 = 15$ non trivial distances which can be obtained from the four operators **Delete**, **Insert**, **Replace** and **Swap**, the SWAP EDIT DISTANCE is the simplest for which no linear time algorithm is known: Wagner and Lowrance [17] described a dynamic program to compute it in time within $O(n^2)$. As the **swap** operator does not remove nor insert any symbol, the SWAP EDIT DISTANCE between two strings of distinct lengths, or of same lengths but with distinct **Parikh vectors**, is always infinite. Such cases can be checked in time within $O(n + m + \sigma)$, and can be ignored as degenerated cases.

The computation of the SWAP EDIT DISTANCE between two strings $S, T \in [1..\sigma]^n$ of same **Parikh vector** $(n_i)_{i \in [1..\sigma]}$ is quite similar to the problem of **SORTING** one string into the other via the exchange of consecutive elements: the **swap** operator is merely reordering the symbols of the strings, and the combined actions of all the **swap** operations can be summarized by a simple permutation over $[1..n]$. Consider the shortest sequence of such **swap** operators transforming S into T , and π the corresponding permutation. The number d of inversions in π , defined as the number of pairs $i, j \in [1..n]$ of positions $i < j$ such that $\pi[j] < \pi[i]$ (i.e. the order of (i, j) is inverse to that of $(\pi[i], \pi[j])$), is exactly the number d of **swaps** required to “reorder” S into T , i.e. the SWAP EDIT DISTANCE d from S to T .

Moffat and Petersson [14] described two sorting algorithms adaptive to the number d of inversions of the input. The first one is the classical sorting algorithm **Insertion Sort**, which sorts an array A with d inversions using only the operator **swaps**, using $O(d) \subseteq O(n^2)$ comparisons and **swaps**. The second one is the adaptive algorithm **Local Insertion Sort**, based on a **Finger Search Tree**, which uses only $O(n(1 + \lg(d/n)))$ comparisons and can be easily modified to *count the number d of inversions*, i.e. the SWAP EDIT DISTANCE d between an array A and its sorted version. We describe below how to take advantage of this algorithm to compute the SWAP EDIT DISTANCE between two arbitrary strings.

First, consider the one-to-one mapping between positions in S and positions in T :

Lemma 3. *Given two strings $S, T \in [1..\sigma]^n$ of same **Parikh vector**, the i -th symbol α in S is mapped to the i -th symbol α in T by the shortest sequence of **Swap** operations transforming S into T .*

Proof. As **Swap** is the only operator available, no symbol is added or removed from S to obtain T , so that there is a one to one mapping between each symbol of S and each symbol of T . Moreover, any sequence of **Swap** operation matching the i -th symbol α in T to the j -th occurrence of α in T can be made shorter if $j \neq i$.

Then, consider how to compute the SWAP EDIT DISTANCE using such a mapping and an index supporting the **rank** and **select** operators:

Theorem 4. *Given two strings $S, T \in [1..\sigma]^n$ of same **Parikh vector** $(n_i)_{i \in [1..\sigma]}$, there is an algorithm computing the SWAP EDIT DISTANCE from S to T via $O(n(1 + \lg(d/n)))$ **rank** and **select** operations.*

Proof. Define the following process to decide if the symbols at positions i and j in S must be inversed during the transformation of S into T minimizing the number of **Swap** operations: Let $a = S[i]$ and $b = S[j]$, so that i and j are the $\text{rank}(S, a, i)$ -th and $\text{rank}(S, b, j)$ -th occurrences of a and b in S , respectively. Let $i' = \text{select}(T, a, \text{rank}(S, a, i))$ and $j' = \text{select}(T, b, \text{rank}(S, b, j))$ be the positions of the corresponding occurrences in T . The symbols will be inversed during the transformation of S into T if and only if (i, j) and (i', j') 's orders are inversed.

“Sorting” S into T using the algorithm **Local Insertion Sort** described by Moffat and Petersson [14] and the process described above to answer comparisons between $\pi(i)$ and $\pi(j)$ yields an algorithm computing the SWAP EDIT DISTANCE using within $O(n(1 + \lg(d/n)))$ **rank** and **select** operations.

This yields as many solutions as there are data structures to support the **rank** and **select** operators, each yielding a distinct computational tradeoff on the previous lemma: we describe two. The first one is based on *inverted posting lists* [21] and an amortized analysis of **doubling search** algorithm [6]:

Corollary 1. *Given two strings $S, T \in [1..\sigma]^n$ of same **Parikh vector** $(n_i)_{i \in [1..\sigma]}$, there is an algorithm computing the SWAP EDIT DISTANCE from S to T in time within $O(n(1 + \lg(d/n)) \lg(n)) \subset O(n^2)$ in the comparison based decision tree model.*

Proof. A simple combination of Theorem 4 with the classical *inverted posting list* implementation [21] of an index supporting the **select** operator in constant time and the **rank** operator via *doubling search* [6].

Note that it should be possible to refine the analysis, as $O(\lg n)$ is a very crude upper bound on the complexity of supporting **rank** or **select**, one should be able to express it in function of n_α , and to amortize it over all **rank** and **select** operations for each symbol α .

The second one is based on the more sophisticated succinct data structure described by Golynski *et al.* [10]:

Corollary 2. *Given two strings $S, T \in [1..\sigma]^n$, there is an algorithm computing the SWAP EDIT DISTANCE from S to T in time within $O(n(1 + \lg(d/n)) \lg \lg \sigma)$.*

Proof. A simple combination of Theorem 4 with the index described by Golynski *et al.* [10] to support the **rank** and **select** operators.

Next, we discuss the minor topic of combining techniques that take advantage of the **Parikh vectors** of the input strings with techniques that take advantage of the output distance.

B Distance Adaptive Computation for all Edit Distances

Abu-Khzam *et al.* [1] described an algorithm computing the INSERT SWAP EDIT DISTANCE d from S to T in time within $O(1.6181^d m)$, which is adaptive to the distance d . We show here that a similar technique can be applied to the other edit distances based on the operators **Delete**, **Insert** and **Replace**, so that to obtain a complexity adaptive to the distance d being computed (Section B.1) and to combine this technique with the ones described previously (Section B.2).

B.1 Distance Adaptive Computation

Lemma 4. *For any edit distance based on a subset of the set of operators $\{\text{Delete}, \text{Insert}, \text{Replace}\}$, there is an algorithm which checks that the distance d from a source string $S \in [1..\sigma]^n$ to a target string $T \in [1..\sigma]^m$ is smaller than a promise $D \in [0..\max\{n, m\}]$ (i.e. if $d \leq D$) in time within $O(D \min\{n, m\})$ and space within $O(n + m)$.*

Theorem 5. *For any edit distance based on a subset of the set of operators $\{\text{Delete}, \text{Insert}, \text{Replace}\}$, there is an algorithm which computes this distance d from a source string $S \in [1..\sigma]^n$ to a target string $T \in [1..\sigma]^m$ in time within $O(d \min\{n, m\})$ and space within $O(n + m)$.*

B.2 Combination with Other Adaptive Techniques

Corollary 3. *There is an algorithm which computes this DELETE INSERT EDIT DISTANCE d from a source string $S \in [1..\sigma]^n$ to a target string $T \in [1..\sigma]^m$ in time within $O(d \min\{n', m'\})$ and space within $O(n + m)$, where $n' = \sum_{\alpha, n_\alpha > 0} n_\alpha$ and $m' = \sum_{\alpha, m_\alpha > 0} m_\alpha$.*