

# Computing the edit distance of a regular language<sup>1</sup>

Stavros Konstantinidis

Department of Mathematics and Computing Science  
Saint Mary's University  
Halifax, Nova Scotia B3H 3C3, Canada  
s.konstantinidis@smu.ca

**Abstract.** The edit distance (or Levenshtein distance) between two words is the smallest number of substitutions, insertions, and deletions of symbols that can be used to transform one of the words into the other. In this paper we consider the problem of computing the edit distance of a regular language (also known as constraint system), that is, the set of words accepted by a given finite automaton. This quantity is the smallest edit distance between any pair of distinct words of the language. We show that the problem is of polynomial time complexity. We distinguish two cases depending on whether the given automaton is deterministic or nondeterministic. In the latter case the time complexity is higher. Incidentally, we also obtain an upper bound on the edit distance of a regular language in terms of the automaton accepting the language.

**Key words:** algorithm, automaton, constraint system, edit distance, Levenshtein distance, regular language.

## 1 Introduction

The problem of measuring the distance, or generally the difference, between words and languages (sets of words) is important in various applications of information processing such as error control in data communications, bio-informatics, speech recognition and spelling correction. The languages of interest are usually regular languages (also called constraint systems), that is, languages defined by finite automata (edge-labeled finite state graphs), but they could be non regular as well. Well-known measures of the difference between two words are the edit (or Levenshtein) distance and the Hamming distance. The edit distance between two words is the smallest number of substitutions, insertions, and deletions of letters required to transform one of the words to the other. Typical problems pertaining to difference measures for words and languages are the following.

1. *The word difference problem:* Compute the edit distance between two given words. The problem can be solved using a dynamic programming algorithm – see [10], for instance.
2. *The error-correction problem:* Given a language of valid (or correct) words, correct a given word to some word of the language that is the least different to the given word. This problem presupposes an agreed measure of word difference such as the edit distance. For (arbitrary) regular languages and for the measure of edit distance the problem was first addressed in [12]. In [7] it is addressed for cases where the language in question is not even regular. In

---

<sup>1</sup>Research supported by grant R220259 of the Natural Sciences and Engineering Research Council of Canada.

[6] and [2] the problem is considered for regular languages and for general word difference measures defined by weighted automata. In coding theory of course, the problem has been addressed extensively for various instances of the word difference measure – this measure is usually implicitly specified by the communications channel that is capable of changing words.

3. *The error-detection capability problem:* Compute the distance (also referred to as self-distance or inner distance) of a given language. This quantity is simply the minimum distance between any pair of distinct words in the language. When the language in question is viewed as a code, the value of the distance represents the maximum number of errors that the code can detect. For example, it is well known that a block code (set of equal length words) can detect up to  $m$  bit substitution errors if and only if the Hamming distance of the code is greater than  $m$ . A similar observation exists for the case of the edit distance [4]. In [3], the authors show how to compute the Hamming distance of a given regular language.

In this work, we address the third problem for the case of the edit distance of a regular language. In this case, the value of the edit distance represents the maximum number  $m$ , say, of substitution, insertion and deletion errors that can be *detected* in the words of the language. This means that no word of the language can be transformed to another different word of the language if up to  $m$  errors are used.

The paper is organized as follows. In the next section we provide the basic notation about words, automata, and edit strings. An edit string is a special word whose symbols are called edit operations. It has been used to define formally the edit distance between ordinary words. Here we also consider languages and finite automata of edit strings as tools for reasoning about the distance of a language. In Section 3, we obtain a few lemmata that will be used to prove the correctness of the main result of the paper. We believe that these lemmata might be of interest in their own right. For example, we show an infinitely-often tight upper bound on the edit distance of a regular language, which depends on the automaton accepting the language. Section 4, contains the main result of the paper, that is, a polynomial time algorithm to compute the edit distance of a given regular language. The time complexity of the algorithm is higher when the language is given via a nondeterministic automaton. Section 5 contains a few concluding remarks and proposes possible directions for future research.

## 2 Basic background and notation

### 2.1 Cardinality, alphabet, word, language

For any set  $S$ , we use the expression  $|S|$  to denote the *cardinality* of the set  $S$ . An *alphabet* is a finite nonempty set  $\Delta$  whose elements are called *symbols*. A *word* or *string* (over  $\Delta$ ) is a finite sequence  $a_1 \cdots a_n$  such that each  $a_i$  is in  $\Delta$ . The *length* of a word  $w$  is denoted by  $|w|$ . The *empty word* is the word of length zero. For any two words  $w_1, w_2$  the expression  $w_1w_2$  denotes the word that obtains by concatenating  $w_1$  and  $w_2$ . A *language* is a set of words. The language of all words is denoted by  $\Delta^*$ .

#### *Notational convention.*

A word  $w$  of length  $n$  can be viewed as a mapping from the index set  $\{1, \dots, n\}$  into  $\Delta$  such that  $w(i)$  is *the symbol of  $w$  at position  $i$* . If  $i > |w|$ , we agree to write that  $w(i) = \perp$ , where  $\perp$  is a

symbol not in the alphabet  $\Delta$ . With this convention, it follows that any two words  $u$  and  $v$  are different if and only if  $u(i) \neq v(i)$  for some positive integer  $i$ .

## 2.2 Finite automaton, computation, diameter, regular language

A (nondeterministic) *finite automaton* is a quintuple  $A = (\Delta, S, s_0, F, T)$  such that  $\Delta$  is the alphabet,  $S$  is the (finite and nonempty) set of states,  $s_0$  is the start state,  $F$  is the set of final states, and  $T$  is the set of transitions, which we denote with expressions of the form  $paq$  such that  $p, q$  are states and  $a$  is a symbol in  $\Delta$ . The automaton  $A$  is said to be deterministic if, for any two transitions of the form  $paq_1, paq_2$ , it is the case that  $q_1 = q_2$ . A *computation* of  $A$  is denoted by an expression of the form  $p_0a_1p_1 \cdots a_np_n$ , for some  $n \geq 0$ , where each  $p_{i-1}a_ip_i$  is a transition of  $A$ . When  $A$  is viewed as a directed edge-labeled graph, a computation is simply a labeled path of  $A$ . The *diameter* of  $A$ , denoted by  $\text{diam}(A)$ , is the largest number of states in a computation  $p_0a_1p_1 \cdots a_np_n$  for which  $p_0$  is the start state and no state occurs more than once, that is,  $i \neq j$  implies  $p_i \neq p_j$ . Obviously the following inequalities about  $\text{diam}(A)$  hold:

$$1 \leq \text{diam}(A) \leq |S|.$$

If, in the computation  $p_0a_1p_1 \cdots a_np_n$ , the state  $p_0$  is the start state and  $p_n$  is a final state then the word  $a_1 \cdots a_n$  is accepted by  $A$ . The *language accepted* by the automaton  $A$  is denoted by  $L(A)$ . The *size* of  $A$ , denoted with  $|A|$ , is the quantity  $|S| + |T|$ , which is the number of states plus the number of transitions. A language  $L$  is called *regular*, or a *constraint system*, if there is a finite automaton accepting  $L$  [5, 13]. If  $A$  and  $A'$  are two finite automata then the expression

$$A \cap A'$$

denotes the finite automaton that obtains using the standard cross product construction such that  $L(A \cap A') = L(A) \cap L(A')$  – see [13]. Moreover,  $|A \cap A'| = O(|A||A'|)$ .

## 2.3 Ordinary word, edit string, weight, edit distance

In this paper we shall use a fixed, but arbitrary, alphabet  $\Sigma$  of *ordinary* symbols, and the alphabet  $E$  of the (*basic*) *edit operations* that depends on  $\Sigma$ . The *empty word* over  $\Sigma$  is denoted by  $\lambda$ , that is,  $\lambda w = w\lambda = w$ , for all words  $w$ . The alphabet  $E$  consists of all symbols  $(x/y)$  such that  $x, y \in \Sigma \cup \{\lambda\}$  and at least one of  $x$  and  $y$  is in  $\Sigma$ . If  $(x/y)$  is in  $E$  and  $x$  is not equal to  $y$  then we call  $(x/y)$  an *error* [1]. We write  $(\lambda/\lambda)$  for the empty word over the alphabet  $E$ . We note that  $\lambda$  is used as a formal symbol in the elements of  $E$ . For example, if  $a$  and  $b$  are in  $\Sigma$  then  $(a/\lambda)(a/b) \neq (a/a)(\lambda/b)$ . The elements of  $E^*$  are called *edit strings* [1], or alignments [6]. The *input* and *output* parts of an edit string  $h = (x_1/y_1) \cdots (x_n/y_n)$  are the words (over  $\Sigma$ )  $x_1 \cdots x_n$  and  $y_1 \cdots y_n$ , respectively. We write  $\text{inp}(h)$  for the input part and  $\text{out}(h)$  for the output part of  $h$ . The expression  $\text{weight}(h)$  denotes the number of errors in  $h$ .

The *edit (or Levenshtein) distance* between two words  $u$  and  $v$ , denoted by  $\text{dist}(u, v)$ , is the smallest number of errors (substitutions, insertions and deletions of symbols) that can transform  $u$  to  $v$ . More formally,

$$\text{dist}(u, v) = \min\{\text{weight}(h) \mid h \in E^*, \text{inp}(h) = u, \text{out}(h) = v\}.$$

For example, for  $\Sigma = \{a, b\}$ , we have that  $\text{dist}(ababa, babbb) = 3$  and the edit string

$$h = (a/\lambda)(b/b)(a/a)(b/b)(a/b)(\lambda/b)$$

is a minimum weight edit string that transforms  $ababa$  to  $babbb$ . In words,  $h$  says that we can use the deletion  $(a/\lambda)$ , the substitution  $(a/b)$ , and the insertion  $(\lambda/b)$  to transform  $ababa$  to  $babbb$ . We note the following fact from [4], for all words  $x, u, v, y$ :

$$\text{dist}(xuy, xvy) = \text{dist}(u, v).$$

If  $L$  is a language containing at least two words then the edit distance of  $L$  is

$$\text{dist}(L) = \min\{\text{dist}(u, v) \mid u, v \in L \text{ and } u \neq v\}.$$

**Definition 1** *Let  $L$  be any set of at least two words. We say that an edit string  $h$  realizes the edit distance of  $L$  if  $\text{weight}(h) = \text{dist}(L)$  and  $\text{inp}(h) \neq \text{out}(h)$  and  $\text{inp}(h), \text{out}(h) \in L$ .*

### 3 Intermediate lemmata

In this section we obtain a few lemmata that will be used to prove the correctness of the main result of the paper. We believe that these lemmata might be of interest in their own right. Moreover the following two paragraphs provide useful technical tools.

*Argument used in proofs.*

In the proofs of the various statements we shall use the following argument about an automaton  $A$  with  $s$  states. If  $P = p_0a_1p_1 \cdots a_kp_k$  is a computation of  $A$  and the number of states  $k+1$  appearing in  $P$  exceeds  $s$ , or  $\text{diam}(A)$  when  $p_0$  is the start state, then at least two states  $p_i$  and  $p_j$ , with  $i < j$ , in  $P$  must be equal. Then we can define a shorter computation  $P'$  of  $A$  that also starts with  $p_0$  and ends with  $p_k$ , by removing the part  $p_i a_i \cdots p_{j-1} a_{j-1}$  of  $P$ . Moreover, for the words  $w_1 = a_1 \cdots a_{i-1}$ ,  $w = a_i \cdots a_j$  and  $w_2 = a_{j+1} \cdots a_k$  we have that  $w$  is nonempty and both of the words  $w_1 w w_2$  and  $w_1 w_2$  are formed in computations of  $A$ . Of course, this argument is well-known, but we include it here for completeness.

*The automaton  $A \cap_E A$ .*

Given two automata  $A, A'$  and a subset  $D$  of  $E$ , the finite automaton  $A \cap_D A'$  [2] accepts all edit strings  $h$  that transform a word of  $L(A)$  into a word of  $L(A')$  using only the edit operations in  $D$ . Moreover, we note that  $|A \cap_D A'| = O(|A||A'|)$ . In our considerations in particular, we shall use the finite automaton  $A \cap_E A$  such that

$$L(A \cap_E A) = \{h \in E^* \mid \text{inp}(h), \text{out}(h) \in L(A)\}.$$

The states of  $A \cap_E A$  are pairs  $(p, q)$ , where  $p$  and  $q$  are states of  $A$ . Hence, if  $A$  has  $s$  states then  $A \cap_E A$  has at most  $s^2$  states.

**Lemma 1** *Let  $A$  be a deterministic finite automaton accepting at least two words. There are distinct words  $u, v$  in  $L(A)$  and an index  $j \leq \text{diam}(A)$  such that  $u(j) \neq v(j)$  and  $\text{dist}(L(A)) = \text{dist}(u, v)$ .*

*Proof.* Suppose that  $\text{dist}(L(A)) = \text{dist}(u, v)$  for some distinct words  $u, v$  in  $L(A)$ . We choose the pair  $(u, v)$  to be minimal as follows: if  $(u_1, v_1)$  is another such pair then  $|u| + |v| \leq |u_1| + |v_1|$ . Let  $j$  be the smallest index for which  $u(j) \neq v(j)$  and assume that  $j > \text{diam}(A)$ . We show that this assumption leads to a contradiction. First note that  $j = |z| + 1$ , where the word  $z$  is the common prefix of length  $j - 1$  of  $u$  and  $v$ . Let  $u'$  and  $v'$  be the words defined by  $u = zu'$  and  $v = zv'$ . Then  $u' \neq v'$  and  $\text{dist}(u, v) = \text{dist}(u', v')$ . As  $|z| \geq \text{diam}(A)$ , the path of  $A$  on which the word  $z$  is formed, which starts from the start state, contains more than  $\text{diam}(A)$  states. This implies that the path contains a repeated state and, therefore, there are words  $w_1, w, w_2$  such that  $|w| > 0$  and  $z = w_1 w w_2$  and the words  $w_1 w_2 u'$  and  $w_1 w_2 v'$  are distinct and in  $L(A)$ . Moreover,  $\text{dist}(w_1 w_2 u', w_1 w_2 v') = \text{dist}(u', v') = \text{dist}(L(A))$ . This, however, contradicts the choice of  $(u, v)$  being minimal. Hence,  $j \leq \text{diam}(A)$ .  $\square$

Unfortunately, the above proof cannot be used if the automaton  $A$  is nondeterministic. This is because there can be two different paths of  $A$  on which the prefix  $z$  is formed and, therefore, the pairs of repeated states in these paths are not identical, in general. Thus, the fact that  $w_1 w w_2 u'$  and  $w_1 w w_2 v'$  are in  $L(A)$  does not imply that also  $w_1 w_2 u'$  and  $w_1 w_2 v'$  are in  $L(A)$ .

**Lemma 2** *Let  $A$  be a (nondeterministic) finite automaton accepting at least two words. There are distinct words  $u, v$  in  $L(A)$  and an index  $j \leq s^2$  such that  $u(j) \neq v(j)$  and  $\text{dist}(L(A)) = \text{dist}(u, v)$ , where  $s$  is the number of states in  $A$ .*

*Proof.* For the sake of contradiction, assume that for all pairs of distinct words  $u, v$  in  $L(A)$  with  $\text{dist}(u, v) = \text{dist}(L(A))$  it is the case that, for any index  $j$ ,  $u(j) \neq v(j)$  implies  $j > s^2$ . For any edit string  $f$  for which the input and output parts are different, we denote with  $i_f$  the smallest value of an index  $i$  for which  $(\text{inp}(f))(i) \neq (\text{out}(f))(i)$ . Now let  $h$  be an edit string that realizes the edit distance of  $L(A)$ . Moreover,  $h$  is chosen to be of minimal length  $|h|$ . By the assumption, it is the case that  $i_h > s^2$ . There are distinct words  $x$  and  $y$ , and alphabet symbols  $a_1, \dots, a_{i_h-1}$  such that

$$\text{inp}(h) = a_1 \cdots a_{i_h-1} x \quad \text{and} \quad \text{out}(h) = a_1 \cdots a_{i_h-1} y$$

and  $\text{dist}(L(A)) = \text{dist}(x, y)$  and  $x(1) \neq y(1)$ . Let  $g$  be an edit string that realizes the edit distance of  $\{x, y\}$ . Define the edit string

$$h' = (a_1/a_1) \cdots (a_{i_h-1}/a_{i_h-1}) g.$$

One can verify that  $h'$  realizes the edit distance of  $L(A)$ . Now the prefix  $(a_1/a_1) \cdots (a_{i_h-1}/a_{i_h-1})$  of  $h'$  appears in some computation, call it  $P$ , of  $A \cap_E A$ . This computation involves  $i_h$  states and, as  $i_h > s^2$ , there is a repeated state in  $P$ . This implies that there is a shorter computation, say  $P'$ , from the first to the last state of  $P$  and, therefore, there is an edit string  $h''$  that is shorter than  $h$  and realizes the edit distance of  $L(A)$ . But this contradicts the assumption about the minimality of  $h$ .  $\square$

Thus, when the automaton  $A$  is nondeterministic the upper bound on the index  $j$  is larger. In the next section, in the context of discussing the complexity of the algorithm for computing the edit distance, we demonstrate that the two upper bounds on the index  $j$  cannot be improved asymptotically, at least.

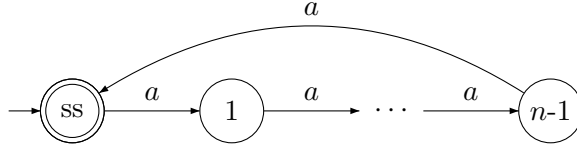


Figure 1: The finite automaton  $C_n$ . The start state is  $ss$ . In pictures for automata, final states are indicated with double lines.

In the next lemma we show that the edit distance of a regular language is upper bounded by the diameter of any automaton accepting the language. This upper bound cannot be improved, in general. For example, in Figure 3, we show a sequence of automata  $(C_n)$  such that  $\text{diam}(C_{n-1}) < \text{diam}(C_n)$  and  $L(C_n)$  consists of the words  $\lambda, a^n, a^{2n}, \dots$ . Moreover,  $\text{dist}(L(C_n)) = \text{diam}(C_n) = n$ .

**Lemma 3** *For any finite automaton  $A$  accepting at least two words,  $\text{dist}(L(A)) \leq \text{diam}(A)$ .*

*Proof.* Let  $(u_1, u_2, \dots)$  be an ordering of  $L(A)$  such that  $|u_i| \leq |u_{i+1}|$  for all  $i$ . The word  $u_2$  can be written as  $u'_1 z$ , for some words  $u'_1$  and  $z$  with  $|u'_1| = |u_1|$ . Obviously,

$$\text{dist}(u_1, u'_1 z) \leq \text{dist}(u_1, u'_1) + \text{dist}(u'_1, u'_1 z) \leq |u'_1| + |z|.$$

Hence,  $\text{dist}(L(A)) \leq |u_2|$ . If  $|u_2| < \text{diam}(A)$ , the statement is true. So assume that  $|u_2| \geq \text{diam}(A)$ . Then there is an accepting computation of  $A$  for  $u_2$  of the form

$$P = p_0 a_1 p_1 \cdots a_i p_i \cdots p_{k-1} a_k p_k \cdots a_n p_n$$

such that  $n+1 > \text{diam}(A)$  and  $k$  is the first index for which  $p_k = p_i$ , for some  $i < k$ . Then there are words  $w_1, w, w_2$  such that  $w$  is nonempty,  $|w_1 w| = k$ ,  $w_1 w w_2 = u_2$  and  $w_1 w_2$  is in  $L(A)$  (in fact  $w_1 w_2$  must be equal to  $u_1$ ). Moreover, it follows that  $\text{dist}(L(A)) \leq |w|$ . Now, as all states  $p_0, \dots, p_{k-1}$  are distinct and  $p_0$  is the start state, we have that  $k \leq \text{diam}(A)$ . This implies  $|w| \leq \text{diam}(A)$  and, therefore,  $\text{dist}(L(A)) \leq \text{diam}(A)$  as required.  $\square$

Next we show the existence of a polynomial size automaton accepting all edit strings  $h$  for which the input part of  $h$  and the output part of  $h$  differ at position  $k$ , for some given positive integer  $k$ . In the expression for the size of this automaton we show explicitly the contribution of the size  $|\Sigma|$  of the alphabet  $\Sigma$  because the alphabet is not related to the given parameter  $k$  in any way.

**Lemma 4** *For any positive integer  $k$ , we can construct a finite automaton  $T_k$  of size  $\Theta(k^2 |\Sigma|^2 + k |\Sigma|^3)$  such that  $T_k$  accepts all edit strings  $h$  for which  $(\text{inp}(h))(k) \neq (\text{out}(h))(k)$ .*

*Proof.* The states of  $T_k$  and their associated meanings are defined as follows.

1.  $[i, j]$  for  $0 \leq i < k$  and  $0 \leq j < k$ : means that the automaton has seen exactly  $i$  input symbols and exactly  $j$  output symbols.
2.  $[ak, j]$  for  $a \in \Sigma$  and  $0 \leq j < k$ : means that the automaton has seen at least  $k$  input symbols and exactly  $j$  output symbols, and the  $k$ -th input symbol was  $a$ .

3.  $[i, ak]$  for  $a \in \Sigma$  and  $0 \leq i < k$ : means that the automaton has seen at least  $k$  output symbols and exactly  $i$  input symbols, and the  $k$ -th output symbol was  $a$ .
4.  $[k, k]$ : means that the automaton has seen at least  $k$  input and at least  $k$  output symbols and the  $k$ -th input symbol was different from the  $k$ -th output symbol.

The start state is  $[0, 0]$  and the final states are  $[ak, j]$ ,  $[i, ak]$ , and  $[k, k]$  for all  $i, j < k$  and  $a \in \Sigma$ . Of course  $T_k$  accepts an edit string  $h$  if and only if  $T_k$  starts at state  $[0, 0]$  and ends its computation on  $h$  at one of the final states. This implies that  $h$  is accepted if and only if  $(\text{inp}(h))(k) \neq (\text{out}(h))(k)$ , as required. The correctness of the construction is established by defining the transitions of  $T_k$  such that the meaning of the states is preserved. This is done next. It is assumed that  $0 \leq i, j \leq k - 1$  and  $a, b, c \in \Sigma$ .

1.  $[i, j](\lambda/a)[i, j + 1]$ , if  $j < k - 1$ . This defines  $k(k - 1)|\Sigma|$  transitions.
2.  $[i, j](\lambda/a)[i, ak]$ , if  $j = k - 1$ . This defines  $k|\Sigma|$  transitions.
3.  $[i, j](a/\lambda)[i + 1, j]$ , if  $i < k - 1$ . This defines  $k(k - 1)|\Sigma|$  transitions.
4.  $[i, j](a/\lambda)[ak, j]$ , if  $i = k - 1$ . This defines  $k|\Sigma|$  transitions.
5.  $[i, j](a/b)[i + 1, j + 1]$ , if  $i, j < k - 1$ . This defines  $(k - 1)^2|\Sigma|^2$  transitions.
6.  $[i, j](a/b)[ak, j + 1]$ , if  $i = k - 1$  and  $j < k - 1$ . This defines  $(k - 1)|\Sigma|^2$  transitions.
7.  $[i, j](a/b)[i + 1, bk]$ , if  $i < k - 1$  and  $j = k - 1$ . This defines  $(k - 1)|\Sigma|^2$  transitions.
8.  $[i, j](a/b)[k, k]$ , if  $i = j = k - 1$  and  $a \neq b$ . This defines  $|\Sigma|^2 - |\Sigma|$  transitions.
9.  $[k, k](x/y)[k, k]$  with  $x, y \in \Sigma \cup \{\lambda\}$  and  $xy \neq \lambda$ . This defines  $|\Sigma|^2 + 2|\Sigma|$  transitions.
10.  $[ak, j](b/\lambda)[ak, j]$ . This defines  $k|\Sigma|^2$  transitions.
11.  $[ak, j](b/c)[ak, j + 1]$ , if  $j < k - 1$ . This defines  $(k - 1)|\Sigma|^3$  transitions.
12.  $[ak, j](b/c)[k, k]$ , if  $j = k - 1$  and  $a \neq c$ . This defines  $|\Sigma|^2(|\Sigma| - 1)$  transitions.
13.  $[ak, j](\lambda/b)[ak, j + 1]$ , if  $j < k - 1$ . This defines  $(k - 1)|\Sigma|^2$  transitions.
14.  $[ak, j](\lambda/b)[k, k]$ , if  $j = k - 1$  and  $a \neq b$ . This defines  $|\Sigma|^2(|\Sigma| - 1)$  transitions.
15.  $[i, ak](\lambda/b)[i, ak]$ . This defines  $k|\Sigma|^2$  transitions.
16.  $[i, ak](b/c)[i + 1, ak]$ , if  $i < k - 1$ . This defines  $(k - 1)|\Sigma|^3$  transitions.
17.  $[i, ak](b/c)[k, k]$ , if  $i = k - 1$  and  $a \neq b$ . This defines  $|\Sigma|^2(|\Sigma| - 1)$  transitions.
18.  $[i, ak](b/\lambda)[i + 1, ak]$ , if  $i < k - 1$ . This defines  $(k - 1)|\Sigma|^2$  transitions.
19.  $[i, ak](b/\lambda)[k, k]$ , if  $i = k - 1$  and  $a \neq b$ . This defines  $|\Sigma|^2(|\Sigma| - 1)$  transitions.

The claim about the size of  $T_k$  follows if one calculates the total number of transitions – this number exceeds the number of states of  $T_k$ .  $\square$

## 4 Computing the distance in polynomial time

In this section we present the main result of the paper, that is, a polynomial time algorithm to compute the edit distance of a given regular language. The time complexity of the algorithm is higher when the language is given via a nondeterministic automaton.

In the first place one could approach the problem as follows. Let  $A$  be the given finite automaton, and let  $B$  be the automaton  $A \cap_E A$  (see Section 3) accepting all edit strings for which the input and output parts are in  $L(A)$ . If we are able to construct an automaton  $T$  accepting all edit strings for which the input and output parts are distinct, then we can compute the automaton  $B \cap T$  that accepts all edit strings for which the input and output parts are distinct and in  $L(A)$ . When  $B \cap T$  is treated as a weighted graph such that the weight of a transition is either 0 or 1, and it is 1 when the label of the transition is an error, then the edit distance of  $L(A)$  is simply the weight of the shortest path in  $B \cap T$ . The question here is whether the required automaton  $T$  exists. We believe that  $T$  does not exist. For the sake of completeness, however, we note that one can construct a finite-state transducer  $T$  realizing all pairs of words  $(u, v)$  with  $u \neq v$ , but when viewed as an automaton,  $T$  does not accept *all* possible edit strings  $h$  whose parts are distinct as desired.

Although the idea described in the preceding paragraph fails, we can build on that idea using the results of the previous section and arrive at the desired algorithm. In the algorithm we shall use the following result from [3].

**Lemma 5** *There is a linear time  $\Theta(|G|)$  algorithm that takes as input a weighted graph  $G$  whose weights on the edges are in  $\{0, 1\}$ , and a vertex  $p$  of  $G$ , and computes, for each vertex of  $G$ , the length of the shortest path to the vertex from  $p$ .*

**Theorem 1** *The following problem is computable in polynomial time.*

**Input:**  $A$  (deterministic or nondeterministic) finite automaton  $A$  accepting at least two words.

**Output:** The edit distance  $\text{dist}(L(A))$ .

*Proof.* Let  $q(A)$  be the quantity  $\text{diam}(A)$  or  $s^2$ , depending on whether  $A$  is deterministic or not, where  $s$  is the number of states in  $A$ . We have the following algorithm.

```

Input = some finite automaton  $A$ ;
dist = diam( $A$ );
 $B = A \cap_E A$ ;
 $m = q(A)$ ;
for each  $j = 1, \dots, m$ 
{
   $G = B \cap T_j$ ;
  weight = ShortestPathWeight( $G$ );
  if (weight < dist) dist = weight;
}
Output = dist;

```

The function call  $\text{ShortestPathWeight}(G)$  treats the automaton  $G$  as a weighted graph and computes the shortest path to every state from the start state. The weight of a transition  $p(x/y)q$  of



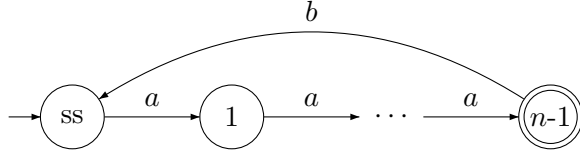


Figure 2: The deterministic finite automaton  $A_n$ .

$G$ , where  $(x/y)$  is an edit operation and  $p, q$  are states, is 1 if  $x \neq y$ , or 0 if  $x = y$ . The algorithm initializes  $\text{dist}$  to the maximum possible value of  $\text{dist}(L(A))$  – see Lemma 3 – and then updates  $\text{dist}$  according to the equation

$$\text{dist}(L(A)) = \min\{\text{weight}(h) \mid h \in L(A \cap_E A) \cap L(T_j), \text{ for some } j \leq q(A)\}.$$

By the lemmata of Section 3, we know that  $\text{dist}(L(A)) = \text{dist}(u, v)$  for some words  $u, v$  in  $L(A)$  with  $u(j) \neq v(j)$  and  $j \leq q(A)$ . Equivalently,  $\text{dist}(L(A))$  is equal to  $\text{weight}(h)$  for some smallest weight edit string  $h$  such that  $h \in L(A \cap_E A)$  and  $(\text{inp}(h))(j) \neq (\text{out}(h))(j)$ , for some  $j \leq q(A)$ . This establishes that the algorithm operates correctly.

We turn now to the time complexity of the algorithm. The function call  $\text{diam}(A)$  that returns the diameter of  $A$  operates in time  $O(|A|)$  using a plain depth first search algorithm that computes, for each state, the length of the shortest (unweighted) path to the state from the start state. The automaton  $B$  can be computed in time  $O(|A|^2)$  and the automaton  $G$  in time  $O(|A|^2|T_j|)$ . By Lemma 5, the function  $\text{ShortestPathWeight}(G)$  operates in time  $O(|A|^2|T_j|)$ . Hence, using also Lemma 4, the dominating term in the time complexity is

$$\sum_{j=1}^m |A|^2(j^2|\Sigma|^2 + j|\Sigma|^3) = O(|A|^2|\Sigma|^2q(|A|)^2(q(A) + |\Sigma|)).$$

□

The question that arises here is whether the quantity  $m = q(A)$ , which is the upper bound of the loop of the algorithm, can be reduced asymptotically. In turn, the question is whether the upper bounds on the index  $j$  that appear in the lemmata of Section 3 can be improved asymptotically. Next we show that this is not possible. Let us establish the following notation for a finite automaton  $A$  accepting at least two words:

- $\text{MinInd}(A)$  denotes the smallest value of an index  $j$  for which there are words  $u, v$  in  $L(A)$  such that  $u(j) \neq v(j)$  and  $\text{dist}(u, v) = \text{dist}(L(A))$ .

The upper bound  $\text{diam}(A)$  of Lemma 1 on the index  $\text{MinInd}(A)$  cannot be improved. More specifically there is a sequence  $(A_n)$  of deterministic finite automata such that  $\text{diam}(A_{n-1}) < \text{diam}(A_n)$  and, for all  $n \geq 2$ ,  $\text{MinInd}(A_n) = \text{diam}(A_n) = n$ . The automaton  $A_n$  is shown in Figure 2.

Similarly, the upper bound  $s^2$  of Lemma 2 on the index  $\text{MinInd}(A)$  cannot be improved asymptotically. More specifically, there is a sequence  $(B_n)$  of nondeterministic finite automata such that  $s_{n-1} < s_n$ , where  $s_n$  is the number of states in  $B_n$ , and  $\text{MinInd}(B_n) = \Theta(s_n^2)$ , as  $n \rightarrow \infty$ . The

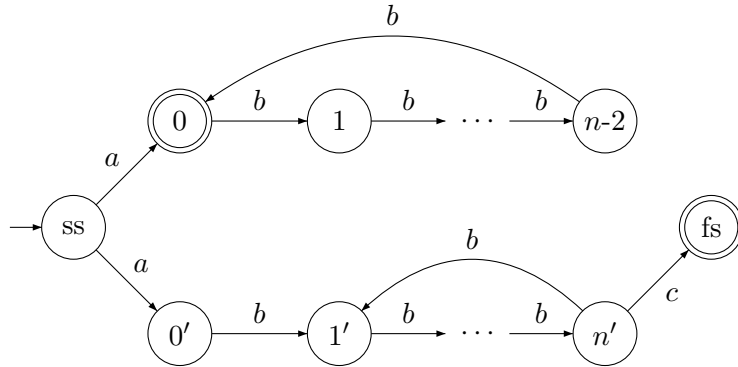


Figure 3: The nondeterministic finite automaton  $B_n$ .

automaton  $B_n$  is shown in Figure 3. It should be clear that  $s_n = 2n + 2$  and the language  $L(B_n)$  consists of all words of the forms  $ab^{i(n-1)}$  and  $ab^{n+jn}c$ , for all  $i, j \geq 0$ . Moreover, for  $n \geq 3$ ,  $\text{dist}(L(B_n)) = 1$  and  $\text{dist}(ab^{i(n-1)}, ab^{(j+1)n}c) = 1$  if and only if  $i(n-1) = (j+1)n$  if and only if  $i = mn$  and  $j+1 = m(n-1)$  for some  $m \geq 1$ . Then, for  $m = 1$ , the words  $ab^{n^2-n}$  and  $ab^{n^2-n}c$  differ at position  $n^2 - n + 2$ , which implies that  $\text{MinInd}(B_n) = n^2 - n + 2 = \Theta(s_n^2)$ .

## 5 Discussion

We have shown why the problem of computing the edit distance of a given regular language is of polynomial time complexity. We have restricted our attention to the case of the unweighted edit distance because this case is closely related to the concept of error detection as discussed in the introduction. However, the methods can be applied even in the case where the edit distance is weighted such that the weights on the errors are positive numbers. In this case, however, the function call  $\text{ShortestPathWeight}(G)$  in the algorithm of Theorem 1 has to invoke Dijkstra's algorithm instead of the algorithm of Lemma 5. This change would increase the time complexity of the algorithm – which of course would remain polynomial.

The first question that arises is whether the time complexity of the algorithms can be improved asymptotically. It appears that this is possible for certain special cases at least. For example, when the given automaton  $A$  is a trellis [11], that is, an automaton with a single final state accepting only words of the same length  $n$ , then  $n = \text{diam}(A) - 1$  and in typical applications  $n$  is much smaller than  $|A|$ .

The next question that arises is whether the results extend to the case of other difference measures for words, in particular, those defined by weighted automata [6, 2]. In this case, the edit strings that one can use to transform words are restricted to only those permitted by the weighted automaton. We note that, for the case of the plain edit distance considered here, two of the arguments that were used in the proofs are (i)  $\text{dist}(xuy, xvy) = \text{dist}(u, v)$ , for all words  $x, u, v, y$ ; and (ii) if  $h$  is any edit string that is used to transform words then also the edit string  $(a_1/a_1) \cdots (a_n/a_n)h$  is permitted for transforming words.

## Acknowledgements

The author is grateful to Nicolae Sântean for providing fruitful comments on this paper.

## References

- [1] L. Kari, S. Konstantinidis. Descriptive complexity of error/edit systems. In: J. Dassow, M. Hoeberechts, H. Jürgensen, D. Wotschke (eds), *Pre-Proceedings of Descriptive Complexity of Formal Systems 2002*, London, Canada, 133–147. To appear in *J. Automata, Languages and Combinatorics* **9** (2004).
- [2] L. Kari, S. Konstantinidis, S. Perron, G. Wozniak, J. Xu, Finite-state error/edit-systems and difference-measures for languages and words. *Tech. Report 2003-01, Dept. Math. and Computing Sci., Saint Mary's University, Canada*, 2003, pp 10. Available electronically at <http://www.stmarys.ca/academic/science/compsci/>
- [3] L. Kari, S. Konstantinidis, S. Perron, G. Wozniak, J. Xu. Computing the Hamming distance of a regular language in quadratic time. *WSEAS Transactions on Information Science & Applications* **1** (2004) 445–449.
- [4] V. I. Levenshtein. Binary codes capable of correcting deletions, insertions, and reversals. *Soviet Physics Dokl.* **10** (1966), 707–710.
- [5] B. H. Marcus, R. M. Roth, P. H. Siegel. Constrained systems and coding for recording channels. In [8], pp. 1635–1764.
- [6] Mehryar Mohri. Edit-distance of weighted automata: general definitions and algorithms. *International Journal of Foundations of Computer Science* **14**(6) (2003), 957-982.
- [7] G. Pighizzini. How hard is computing the edit distance? *Information and Computation* **165** (2001), 1–13.
- [8] V. S. Pless, W. C. Huffman (eds). *Handbook of Coding Theory*, Elsevier, 1998.
- [9] G. Rozenberg, A. Salomaa (eds). *Handbook of Formal Languages, vol. I*. Springer, Berlin, 1997.
- [10] D. Sankoff, J. Kruskal (eds). *Time Warps, String Edits, and Macromolecules: The Theory and Practice of Sequence Comparison*. CSLI Publications, 1999.
- [11] A. Vardy. Trellis structure of codes. In [8], pp. 1989–2117.
- [12] R. A. Wagner. Order- $n$  correction for regular languages. *Communications of the ACM* **17**(5) (1974), 265–268.
- [13] S. Yu. Regular languages. In [9], pp 41–110.