

Bit-parallel approach to approximate string matching in compressed texts

Matsumoto, Tetsuya
Department of Informatics, Kyushu University

Kida, Takuya
Department of Informatics, Kyushu University

Takeda, Masayuki
Department of Informatics, Kyushu University

Shinohara, Ayumi
Department of Informatics, Kyushu University

他

<http://hdl.handle.net/2324/3032>

出版情報 : DOI Technical Report. 174, 2000-04. 九州大学大学院システム情報科学研究所情報理学部門
バージョン :
権利関係 :



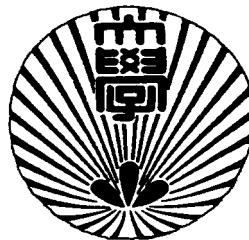
DOI Technical Report

Bit-parallel approach to approximate string matching in compressed texts

by

TETSUYA MATSUMOTO, TAKUYA KIDA, MASAYUKI
TAKEDA, AYUMI SHINOHARA AND SETSUO ARIKAWA

April 2000



Department of Informatics
Kyushu University
Fukuoka 812-8581, Japan

Email: tetsuya@i.kyushu-u.ac.jp Phone: +81-92-642-2697

Bit-parallel approach to approximate string matching in compressed texts

Tetsuya Matsumoto Takuya Kida Masayuki Takeda
Ayumi Shinohara Setsuo Arikawa

Department of Informatics, Kyushu University 33, Fukuoka 812-8581, Japan
{ tetsuya, kida, takeda, ayumi, arikawa } @i.kyushu-u.ac.jp

Abstract

In this paper, we address the problem of approximate string matching on compressed text. We consider this problem for a text string described in terms of *collage system*, which is a formal system proposed by Kida et al. (1999) that captures various dictionary-based compression methods. We present an algorithm that exploits bit-parallelism, assuming that our problem fits in a single machine word, e.g., $(m - k + 1)(k + 1) \leq L$, where m is the pattern length, k is the number of allowed errors, and L is the length in bits of the machine word. For a class of *simple* collage systems, the algorithm runs in $O(k^2(\|\mathcal{D}\| + |\mathcal{S}|) + km)$ time using $O(k^2\|\mathcal{D}\|)$ space, where $\|\mathcal{D}\|$ is the size of dictionary \mathcal{D} and $|\mathcal{S}|$ is the number of tokens in \mathcal{S} . The LZ78 and the LZW compression methods are covered by this class. Since we can regard $n = \|\mathcal{D}\| + |\mathcal{S}|$ as the compressed length, the time and the space complexities are $O(k^2n + km)$ and $O(k^2n)$, respectively.

1 Introduction

The problem of approximate string matching is to find the locations of approximate occurrences P' of a pattern P in a text T such that the edit distance between P and P' is $\leq k$. This problem has been studied extensively. See an excellent survey paper by Navarro [7].

Baeza-Yates and Navarro [2] presented a very fast algorithm for on-line approximate string matching, which is based on the simulation of a nondeterministic finite automaton (NFA) built from the pattern P and the number k of allowed errors. Figure 1 shows the NFA for the pattern $P = \text{POKEMON}$ with $k = 4$. The simulation algorithm exploits bit-parallelism. The idea is essentially similar to the work of Wu and Manber [11], but the search time is improved by packing the automaton states differently. That is, the algorithm simulates NFA using diagonals instead of

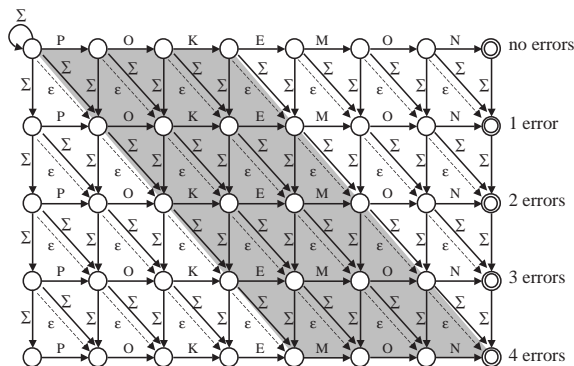


Figure 1: NFA for approximate string matching for the pattern POKEMON with $k = 4$.

rows. Since there are $m - k + 1$ complete diagonals and the others are not really necessary, we need only $(m - k)(k + 2)$ bits for a pattern of length m . Thus the automaton states can be packed into a machine word for a short pattern, and the algorithm achieves an $O(n)$ time search. Recall that the Wu-Manber algorithm uses $k + 1$ machine words to encode the automaton states, and requires $O(kn)$ time.

On the other hand, the compressed pattern matching attracts recently special concern where the goal is to perform the exact string matching in a compressed text without decompressing it. The problem has been studied in 1990's by several researchers mainly for dictionary-based compression methods, such as the Ziv-Lempel family (e.g., LZ77 [12], LZ78 [13], LZW [10]). In [5], we introduced *collage systems*, a formal system to represent a text string, that captures various dictionary-based compressions. Within this framework, we generalized the existing compressed pattern matching algorithms and unified the concepts into a general algorithm. Thus *any* compression method covered by the framework has now an compressed pattern matching algorithm as an instance.

However, the studies on the compressed pattern matching have been undertaken only for the exact pattern matching. The approximate string matching on compressed text has been an open problem since advocated in [1]. The first solution to this problem was very recently given by Kärkkäinen et al. [4]. They addressed the LZ78 and the LZW compression methods. The proposed algorithm runs in $O(mkn + r)$ time, where n is the compressed text length, m is the pattern length, k is the number of allowed errors, and r is the number of pattern occurrences with k errors or less. For the existence problem it needs $O(mkn)$ time and space. The algorithm is based on the simulation of the classical dynamic programming algorithm for approximate string matching [9]. The experimental results reported in [4]

show that, for small k , the algorithm runs faster than the method of decompressing the text on the fly and searching over it using the classical dynamic programming method that takes $O(mN)$ time, where N is the original text length.

In this paper, we take the NFA simulation approach using bit-parallelism, instead of the dynamic programming approach. We present a new algorithm for approximate string matching to solve the existence problem. Like the work of Baeza-Yates and Navarro [2], we assume that our problem fits in a single machine word, i.e., $(m-k+1)(k+1) \leq L$, where L is the length in bits of the machine word. Our algorithm runs on a simple collage system $\langle \mathcal{D}, \mathcal{S} \rangle$ in $O(k^2(\|\mathcal{D}\| + |\mathcal{S}|) + km)$ time using $O(k^2\|\mathcal{D}\|)$ space. It should be emphasized that our algorithm works for *any* compression method that covered by the class of simple collage systems. For example, this class covers the LZ78 and the LZW compression methods. Since we can regard $n = \|\mathcal{D}\| + |\mathcal{S}|$, the time and the space complexities are $O(k^2n + km)$ and $O(k^2n)$, respectively. Thus our algorithm could run faster than the one due to Kärkkäinen et al., especially when k is much smaller than m , although ours cannot cope with the all occurrence problem.

2 Preliminaries

Let Σ be a finite alphabet. An element of Σ^* is called a *string*. Strings x , y , and z are said to be a *prefix*, *factor*, and *suffix* of the string $u = xyz$, respectively. The length of a string u is denoted by $|u|$. The empty string is denoted by ε , that is, $|\varepsilon| = 0$. The i th symbol of a string u is denoted by $u[i]$ for $1 \leq i \leq |u|$, and the factor of a string u that begins at position i and ends at position j is denoted by $u[i : j]$ for $1 \leq i \leq j \leq |u|$. For convenience, let $u[i : j] = \varepsilon$ for $j < i$. Denote by u^R the reversed string of a string u . For a string u and a non-negative integer i , the string obtained by removing the length i prefix (resp. suffix) from u is denoted by ${}^{[i]}u$ (resp. $u^{[i]}$). That is, ${}^{[i]}u = u[i + 1 : |u|]$ and $u^{[i]} = u[1 : |u| - i]$. Denote by $D(x, y)$ the edit distance between two strings x and y .

3 Collage system

In a dictionary-based compression, a text string is described by a pair of a *dictionary* and a sequence of *tokens*, each of which represents a phrase defined in the dictionary. Kida et al. [5] introduced a unifying framework, named collage system, which abstracts various dictionary-based methods, such as the Lempel-Ziv family, the RE-PAIR [6], and static dictionary methods. In [5] they presented a general

compressed pattern matching algorithm for the framework, which is based on the simulation of the KMP automaton. This implies that *any* compression method covered by the framework has a compressed pattern matching algorithm as an instance.

A *collage system* is a pair $\langle \mathcal{D}, \mathcal{S} \rangle$ defined as follows: \mathcal{D} is a sequence of assignments $X_1 = \text{expr}_1; X_2 = \text{expr}_2; \dots; X_\ell = \text{expr}_\ell$, where each X_k is a token (or a variable) and expr_k is any of the form

$$\begin{array}{ll} a & \text{for } a \in \Sigma \cup \{\varepsilon\}, \quad (\text{primitive assignment}) \\ X_i X_j & \text{for } i, j < k, \quad (\text{concatenation}) \\ [j]X_i & \text{for } i < k \text{ and an integer } j, \quad (\text{prefix truncation}) \\ X_i^{[j]} & \text{for } i < k \text{ and an integer } j, \quad (\text{suffix truncation}) \\ (X_i)^j & \text{for } i < k \text{ and an integer } j. \quad (j \text{ times repetition}) \end{array}$$

Each token represents a string obtained by evaluating the expression as it implies. The strings represented by tokens are called *phrases*. As we want to distinguish a token from the phrase it represents, we denote by $X.u$ the phrase represented by a token X . The *size* of \mathcal{D} is the number ℓ of assignments and denoted by $\|\mathcal{D}\|$. Define the *height* of a token X to be the height of the syntax tree whose root is X . The *height* of \mathcal{D} is defined by $\text{height}(\mathcal{D}) = \max\{\text{height}(X) \mid X \text{ in } \mathcal{D}\}$. It expresses the maximum dependency of the tokens in \mathcal{D} .

On the other hand, $\mathcal{S} = X_{i_1}, X_{i_2}, \dots, X_{i_n}$ is a sequence of tokens defined in \mathcal{D} . We denote by $|\mathcal{S}|$ the number n of tokens in \mathcal{S} . The collage system represents a string obtained by concatenating the phrases represented by $X_{i_1}, X_{i_2}, \dots, X_{i_n}$.

According to the framework, text strings compressed by LZW and RE-PAIR can be represented as follows.

LZW. $\mathcal{S} = X_{i_1}, X_{i_2}, \dots, X_{i_n}$ and \mathcal{D} is as follows:

$$\begin{aligned} X_1 &= a_1; & X_2 &= a_2; & \dots; & X_q &= a_q; \\ X_{q+1} &= X_{i_1} X_{\sigma(i_2)}; & X_{q+2} &= X_{i_2} X_{\sigma(i_3)}; & \dots; \\ X_{q+n-1} &= X_{i_{n-1}} X_{\sigma(i_n)}, \end{aligned}$$

where the alphabet is $\Sigma = \{a_1, \dots, a_q\}$, $1 \leq i_1 \leq q$, and $\sigma(\ell)$ denotes the integer k , $1 \leq k \leq q$, such that a_k is the first symbol of the phrase $X_\ell.u$. \mathcal{S} is encoded as a sequence of integers i_1, i_2, \dots, i_n in which an integer i_j is represented in $\lceil \log_2(q+j) \rceil$ bits, while \mathcal{D} is not encoded since it can be obtained from \mathcal{S} .

RE-PAIR. $\mathcal{S} = X_{i_1}, X_{i_2}, \dots, X_{i_n}$, and \mathcal{D} is as follows:

$$\begin{aligned} X_1 &= a_1; & X_2 &= a_2; & \dots; & X_q &= a_q; \\ X_{q+1} &= X_{\ell(1)} X_{r(1)}; & X_{q+2} &= X_{\ell(2)} X_{r(2)}; & \dots; \\ X_{q+s} &= X_{\ell(s)} X_{r(s)}, \end{aligned}$$

where $\Sigma = \{a_1, \dots, a_q\}$. \mathcal{D} and \mathcal{S} are encoded using some appropriate encoding.

Theorem 1 (Kida et al. [5]) *The problem of pattern matching for a text $\langle \mathcal{D}, \mathcal{S} \rangle$ can be solved in $O(\text{height}(\mathcal{D})(\|\mathcal{D}\| + |\mathcal{S}|) + m^2 + r)$ time using $O(\|\mathcal{D}\| + m^2)$ space, where r is the number of pattern occurrences. The factor $\text{height}(\mathcal{D})$ can be dropped if \mathcal{D} contains no truncation.*

Since the factor $\|\mathcal{D}\| + |\mathcal{S}|$ is considered to be the compressed text length, the algorithm runs in linear time proportional to the compressed text length if \mathcal{D} contains no truncation. The result coincides with the observation by Navarro and Raffinot [8] that LZ77 is not suitable for compressed pattern matching compared with LZ78 compression.

The following definitions will be needed for Sections 4 and 5.

Definition 1 *A collage system is said to be regular if it contains neither repetition nor truncation. A regular collage system is said to be simple if, for every assignment $X = YZ$, $|Y.u| = 1$ or $|Z.u| = 1$.*

Note that the collage systems for the RE-PAIR compression are regular, and those for the LZW/LZ78 compressions are simple.

4 Basic idea

We give an overview of our algorithm for approximate string matching, which is based on the simulation of the NFA. First we give a formal description of the NFA for approximate string matching of pattern $P = P[1 : m]$ with k errors. Denote by \mathbf{N} the set of natural numbers $0, 1, \dots$. For any $S \subseteq \mathbf{N} \times \mathbf{N}$ and any $\langle x, y \rangle \in \mathbf{N} \times \mathbf{N}$, let $S \oplus \langle x, y \rangle = \{\langle i + x, j + y \rangle \mid \langle i, j \rangle \in S\}$. Let

$$Q = \{\langle i, j \rangle \mid 0 \leq i \leq m, 0 \leq j \leq k\}.$$

Define the mapping $\delta^{\text{ed}} : Q \times (\Sigma \cup \{\varepsilon\}) \rightarrow 2^Q$ by

$$\begin{aligned} \delta^{\text{ed}}(\langle i, j \rangle, \varepsilon) &= \{\langle i + 1, j + 1 \rangle \mid i + 1 \leq m, j + 1 \leq k\}, \\ \delta^{\text{ed}}(\langle i, j \rangle, a) &= \{\langle i + 1, j \rangle \mid i + 1 \leq m, P[i + 1] = a\} \\ &\quad \cup \{\langle i, j + 1 \rangle \mid j + 1 \leq k\} \\ &\quad \cup \{\langle i + 1, j + 1 \rangle \mid i + 1 \leq m, j + 1 \leq k\}, \end{aligned}$$

where $\langle i, j \rangle \in Q$ and $a \in \Sigma$. Define the mapping $\delta^{\text{pm}} : Q \times (\Sigma \cup \{\varepsilon\}) \rightarrow 2^Q$ by

$$\begin{aligned} \delta^{\text{pm}}(\langle i, j \rangle, \varepsilon) &= \delta^{\text{ed}}(\langle i, j \rangle, \varepsilon), \\ \delta^{\text{pm}}(\langle i, j \rangle, a) &= \delta^{\text{ed}}(\langle i, j \rangle, a) \cup \{\langle 0, 0 \rangle \mid i = j = 0\}, \end{aligned}$$

where $\langle i, j \rangle \in Q$ and $a \in \Sigma$. Let $F = \{\langle m, j \rangle \mid 0 \leq j \leq k\}$.

Proposition 1 *The NFA specified by $(Q, \Sigma, \delta^{\text{ed}}, \{(0, 0)\}, F)$ accepts the language $\{w \in \Sigma^* \mid D(P, w) \leq k\}$.*

Proposition 2 *The NFA specified by $(Q, \Sigma, \delta^{\text{pm}}, \{(0, 0)\}, F)$ accepts the language $\Sigma^*\{w \in \Sigma^* \mid D(P, w) \leq k\}$.*

Let $\hat{\delta}^{\text{ed}}$ and $\hat{\delta}^{\text{pm}}$, respectively, be the functions δ^{ed} and δ^{pm} extended over the domain $2^Q \times \Sigma^*$ in the standard way (see Section 2.4 of [3]).

Definition 2 *Define the function $Jump : 2^Q \times \mathcal{D} \rightarrow 2^Q$ by*

$$Jump(S, t) = \hat{\delta}^{\text{pm}}(S, t.u).$$

Define the function $Output : 2^Q \times \mathcal{D} \rightarrow \{true, false\}$ by

$$\begin{aligned} Output(S, t) &= true \\ \Leftrightarrow & \text{there exists a nonempty prefix } w \text{ of } t.u \text{ such that } \hat{\delta}^{\text{pm}}(S, w) \cap F \neq \emptyset. \end{aligned}$$

The basic idea of our algorithm is to simulate the NFA in compressed text, by packing $S \in 2^Q$ in a machine word. We need only $(m - k + 1)(k + 1)$ states out of Q since there are $m - k + 1$ complete diagonals and the others are not necessary as pointed out by Baeza-Yates and Navarro [2]. That is, we need only the states in

$$Q_{\text{core}} = \{\langle i, j \rangle \in Q \mid 0 \leq i \leq m - k + 1, 0 \leq j \leq k\}.$$

Assume that the given pattern $P = P[1 : m]$ is small so that $(m - k + 1)(k + 1)$ is not greater than the machine word length in bits. Then the union, the intersection, and the \oplus operations on the subsets of Q_{core} can be performed in $O(1)$ time. As shown later, we can build the functions $Jump$ and $Output$ in $O(k^2\|\mathcal{D}\| + km)$ time using $O(k^2\|\mathcal{D}\|)$ space, so that they return their values in $O(k^2)$ time, for a simple collage system $\langle \mathcal{D}, \mathcal{S} \rangle$. Figure 2 gives an overview of the algorithm using the function $Jump$ and $Output$.

We thus have the next result under the assumption that $(m - k + 1)(k + 1) \leq L$, where L is the machine word length in bits.

Theorem 2 *Our algorithm runs in $O(k^2(\|\mathcal{D}\| + |\mathcal{S}|) + km)$ time using $O(k^2\|\mathcal{D}\|)$ space for a simple collage system $\langle \mathcal{D}, \mathcal{S} \rangle$.*

5 Algorithm in detail

In this section, we discuss the construction of the two functions $Jump$ and $Output$. Throughout this section, we assume that the union, the intersection, and the \oplus operations for the subsets of Q can be performed in $O(1)$ time.

```

Input.    A text string represented as a pair of  $\mathcal{D}$  and  $S = S[1 : n]$ ,
            and a pattern  $P[1 : m]$ .
Output.  A boolean value that indicates whether there is an approximate
            occurrence of  $P$  with  $k$  errors or less.
begin
  /* Preprocessing */
  Build Jump and Output from the pattern  $P$  and the dictionary  $\mathcal{D}$ ;

  /* Text scanning */
   $R := I(\varepsilon)$ ;    /*  $I(w)$  is defined in Section 5 */
  for  $\ell := 1$  to  $n$  do begin
    if  $Output(R, S[\ell]) = true$  then return true;
     $R := Jump(R, S[\ell])$ 
  end;
  return false
end.

```

Figure 2: Overview of our algorithm.

5.1 Construction of function *Jump*

Definition 3 For any $w \in \Sigma^*$, let $I(w) = \hat{\delta}^{\text{pm}}(\{\langle 0, 0 \rangle\}, w)$.

We can prove the following relation between the two transition functions $\hat{\delta}^{\text{pm}}$ and $\hat{\delta}^{\text{ed}}$.

Lemma 1 For any $S \in 2^Q$ and any $w \in \Sigma^*$,

$$\hat{\delta}^{\text{pm}}(S, w) = \begin{cases} \hat{\delta}^{\text{ed}}(S, w) \cup I(w), & \text{if } \langle 0, 0 \rangle \in S; \\ \hat{\delta}^{\text{ed}}(S, w), & \text{otherwise.} \end{cases}$$

Note that $I(w)$ corresponds to the existence of the self-loop on the initial state $\langle 0, 0 \rangle$.

Definition 4 For any $a \in \Sigma$, let

$$\begin{aligned} G(a) &= \{\langle i, j \rangle \mid \langle i, j \rangle \in Q, i \neq 0, P[i] = a\}, \\ H(a) &= \{\langle i, j \rangle \mid \langle i, j \rangle \in Q, P[i'] = a \text{ for some } i' \text{ with } 1 \leq i' \leq i\}. \end{aligned}$$

Because we need only the states in Q_{core} and encode each of the sets $G(a) \cap Q_{\text{core}}$ and $H(a) \cap Q_{\text{core}}$ into a machine word, these tables can be built in $O(|\Sigma| + km)$ time using $O(|\Sigma|)$ space in a simple way.

Lemma 2 For any $a \in \Sigma$ and any $u, v \in \Sigma^*$,

$$\begin{aligned} I(a) &= Q \cap \left(E \cup (E \oplus \langle 0, 1 \rangle) \cup ((E \oplus \langle 1, 0 \rangle) \cap H(a)) \right) \\ I(uv) &= \hat{\delta}^{\text{ed}}(I(u), v) \cup I(v), \end{aligned}$$

where $E = \{\langle i, j \rangle \in Q \mid i = j\}$.

The above lemma implies that, for a regular collage system, we can build the table which stores $I(t.u)$ for the tokens t in \mathcal{D} by calling the function $\hat{\delta}^{\text{ed}}$ $O(\|\mathcal{D}\|)$ times. Thus we concentrate on how to compute the function $\hat{\delta}^{\text{ed}}$.

Definition 5 For any $w \in \Sigma^*$ and any integers ℓ and d ($0 \leq \ell \leq m$ and $0 \leq d \leq k$), let

$$V(w; \ell, d) = \left\{ i \mid 0 \leq i \leq m, 0 \leq i - \ell \leq m, \langle i, d \rangle \in \hat{\delta}^{\text{ed}}(\{\langle i - \ell, 0 \rangle\}, w) \right\},$$

and let

$$M(w; \ell, d) = \left\{ \langle i, j \rangle \in Q \mid i \in V(w; \ell, d) \right\}.$$

Then we have the following lemma.

Lemma 3 For any $S \in 2^Q$ and any $w \in \Sigma^*$,

$$\hat{\delta}^{\text{ed}}(S, w) = \bigcup_{\substack{\max(0, |w| - k) \leq \ell \leq \min(|w| + k, m) \\ 0 \leq d \leq k}} (S \oplus \langle \ell, d \rangle) \cap M(w; \ell, d).$$

Proof. Let $\langle i, j \rangle \in \hat{\delta}^{\text{ed}}(S, w)$. There exists $\langle i', j' \rangle \in S$ such that $\langle i, j \rangle \in \hat{\delta}^{\text{ed}}(\{\langle i', j' \rangle\}, w)$. Let $\ell = i - i'$ and $d = j - j'$. It is easy to see that $\langle i, j \rangle \in S \oplus \langle \ell, d \rangle$, and that $0 \leq \ell \leq m$ and $0 \leq d \leq k$. By the definition of $\hat{\delta}^{\text{ed}}$ we see that $\langle i, d \rangle \in \hat{\delta}^{\text{ed}}(\{\langle i', 0 \rangle\}, w)$, and therefore $\langle i, j \rangle \in M(w; \ell, d)$. By Proposition 1, $\langle i, d \rangle \in \hat{\delta}^{\text{ed}}(\{\langle i - \ell, 0 \rangle\}, w)$ implies that $D(P[i - \ell + 1 : i], w) \leq d$. Hence $|w| - d \leq |P[i - \ell + 1 : i]| \leq |w| + d$, and then $|w| - k \leq \ell \leq |w| + k$. Thus we have

$$\hat{\delta}^{\text{ed}}(S, w) \subseteq \bigcup_{\substack{\max(0, |w| - k) \leq \ell \leq \min(|w| + k, m) \\ 0 \leq d \leq k}} (S \oplus \langle \ell, d \rangle) \cap M(w; \ell, d).$$

On the other hand, let $\langle i, j \rangle \in (S \oplus \langle \ell, d \rangle) \cap M(w; \ell, d)$ for some ℓ and d such that $\max(0, |w| - k) \leq \ell \leq \min(|w| + k, m)$ and $0 \leq d \leq k$. We can prove in a similar manner that $\langle i, j \rangle \in \hat{\delta}^{\text{ed}}(S, w)$. \blacksquare

Let \bar{M} be the table which stores $M(t.u; \ell, d)$ for the tokens t in \mathcal{D} , and the integers ℓ, d such that $\max(1, |t.u| - k) \leq \ell \leq \min(m, |t.u| + k)$ and $0 \leq d \leq k$. The size of \bar{M} is $O(k^2 \|\mathcal{D}\|)$. The above lemma implies that, for any $S \in 2^Q$ and any token t , we can obtain the value $\hat{\delta}^{\text{ed}}(S, t.u)$ in $O(k^2)$ time using the table \bar{M} . We now consider the computation of the table \bar{M} .

Lemma 4 *For a simple collage system, we can build the table \bar{M} in $O(k^2\|\mathcal{D}\|+km)$ time using $O(k^2\|\mathcal{D}\|)$ space.*

Proof. We can compute the table by the recurrence

$$M(\varepsilon; \ell, d) = \begin{cases} (Q \oplus \langle \ell, \ell \rangle) \cap Q, & \text{if } \ell = d \\ \emptyset, & \text{otherwise} \end{cases}$$

$$M(ua; \ell, d) = Q \cap \left(M(ua; \ell - 1, d - 1) \oplus \langle 1, 0 \rangle \right. \\ \cup M(u; \ell - 1, d - 1) \oplus \langle 1, 0 \rangle \\ \cup M(u; \ell, d - 1) \\ \left. \cup (M(u; \ell - 1, d) \oplus \langle 1, 0 \rangle) \cap G(a) \right),$$

$$M(au; \ell, d) = M(au; \ell - 1, d - 1) \\ \cup M(u; \ell - 1, d - 1) \\ \cup M(u; \ell, d - 1) \\ \cup (M(u; \ell - 1, d) \cap (G(a) \oplus \langle \ell - 1, 0 \rangle)),$$

where $u \in \Sigma^*$, $a \in \Sigma$, and $\ell \geq 0$ and $d \geq 0$. (We assume, for the sake of convenience, that $M(w; i, j) = \emptyset$ for all $w \in \Sigma^*$ if $i < 0$ or $j < 0$.) \blacksquare

Unfortunately, for a regular collage system, we have not devised a way of computing the table \bar{M} . Thus our result is only for simple collage systems.

Lemma 5 *For a simple collage system, we can build the function *Jump* in $O(k^2\|\mathcal{D}\|+km)$ time using $O(k^2\|\mathcal{D}\|)$ space, so that it returns its value in $O(k^2)$ time.*

5.2 Construction of function *Output*

Definition 6 *For any $S \in 2^{\mathcal{Q}}$ and any token t in \mathcal{D} , define *Output*₁ and *Output*₂ by*

$$\begin{aligned} \text{Output}_1(t) &= \text{true} \\ \Leftrightarrow & \text{the string } t.u \text{ contains an approximate occurrence of } P. \\ \text{Output}_2(S, t) &= \text{true} \\ \Leftrightarrow & \text{there exist an integer } \ell > 0 \text{ and a non-empty prefix } w \text{ of } t.u \text{ such that} \\ & \langle m - \ell, j \rangle \in S \text{ and } j + D(P[m - \ell + 1 : m], w) \leq k. \end{aligned}$$

Definition 7 *For any $S \in 2^{\mathcal{Q}}$ and any token X in \mathcal{D} , define *Occ*^{*} by*

$$\begin{aligned} \text{Occ}^*(X) &= \text{true} \\ \Leftrightarrow & \text{there exist tokens } Y, Z \text{ with } X = YZ \in \mathcal{D} \text{ and } u, v \in \Sigma^* \text{ such that} \\ & u \text{ is a suffix of } Y.u, v \text{ is a prefix of } Z.u, \text{ and } D(uv, P) \leq k. \end{aligned}$$

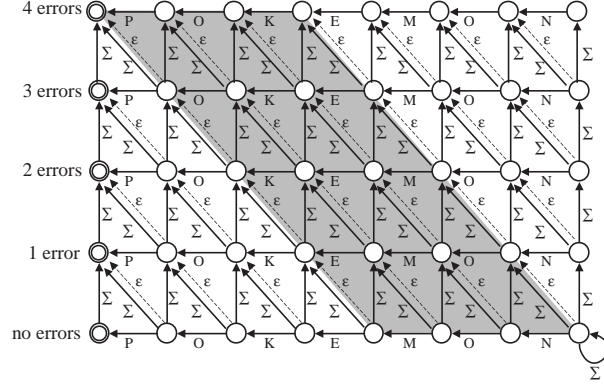


Figure 3: NFA for approximate string matching of the reversed pattern $P^R = \text{NOMEKOP}$ with $k = 4$.

Let us consider the computation of Occ^* . We build the NFA for approximate pattern matching for the reversed pattern P , which plays a key role in the computation of Occ^* . Figure 3 shows the NFA for $P^R = \text{NOMEKOP}$ with $k = 4$. Note that the NFA in Fig. 3 is the same as the one obtained by reversing the direction of the arcs of the NFA in Fig. 1 except for the self-loop on the initial-state. The precise discussion is as follows. Define the mapping $\delta_{\text{rev}}^{\text{ed}} : Q \times (\Sigma \cup \{\varepsilon\}) \rightarrow 2^Q$ by

$$\begin{aligned} \delta_{\text{rev}}^{\text{ed}}(\langle i, j \rangle, \varepsilon) &= \{\langle i-1, j-1 \rangle \mid 0 \leq i-1, 0 \leq j-1\}, \\ \delta_{\text{rev}}^{\text{ed}}(\langle i, j \rangle, a) &= \{\langle i-1, j \rangle \mid 0 \leq i-1, P[i] = a\} \\ &\quad \cup \{\langle i, j-1 \rangle \mid 0 \leq j-1\} \\ &\quad \cup \{\langle i-1, j-1 \rangle \mid 0 \leq i-1, 0 \leq j-1\}, \end{aligned}$$

where $\langle i, j \rangle \in Q$ and $a \in \Sigma$. Define the mapping $\delta_{\text{rev}}^{\text{pm}} : Q \times (\Sigma \cup \{\varepsilon\}) \rightarrow 2^Q$ by

$$\begin{aligned} \delta_{\text{rev}}^{\text{pm}}(\langle i, j \rangle, \varepsilon) &= \delta_{\text{rev}}^{\text{ed}}(\langle i, j \rangle, \varepsilon), \\ \delta_{\text{rev}}^{\text{pm}}(\langle i, j \rangle, a) &= \delta_{\text{rev}}^{\text{ed}}(\langle i, j \rangle, a) \cup \{\langle m, k \rangle \mid i = m, j = k\}, \end{aligned}$$

where $\langle i, j \rangle \in Q$ and $a \in \Sigma$. Let $F_{\text{rev}} = \{\langle 0, j \rangle \mid 0 \leq j \leq k\}$. Obviously we have:

Proposition 3 *The NFA specified by $(Q, \Sigma, \delta_{\text{rev}}^{\text{ed}}, \{\langle m, k \rangle\}, F_{\text{rev}})$ accepts the language $\{w \in \Sigma^* \mid D(P^R, w) \leq k\}$.*

Proposition 4 *The NFA specified by $(Q, \Sigma, \delta_{\text{rev}}^{\text{pm}}, \{\langle m, k \rangle\}, F_{\text{rev}})$ accepts the language $\Sigma^* \{w \in \Sigma^* \mid D(P^R, w) \leq k\}$.*

Lemma 6 *For any strings $u, v \in \Sigma^*$, the string uv contains an approximate occurrence of P with at most k errors if and only if*

$$I(u) \cap I_{\text{rev}}(v^R) \neq \emptyset,$$

where $I_{\text{rev}}(w) = \hat{\delta}_{\text{rev}}^{\text{pm}}(\langle m, k \rangle, w)$ for $w \in \Sigma^*$.

Lemma 7 *For a simple collage system, we can build in $O(k^2\|\mathcal{D}\|+km)$ time using $O(k^2\|\mathcal{D}\|)$ space the table which stores the values $Occ^*(t)$ for the tokens t in \mathcal{D} .*

Proof. All tokens X are defined in \mathcal{D} as either $X = a$ or $X = YZ$, where $a \in \Sigma$ and Y, Z are tokens. When $X = a$, $Occ^*(X) = \emptyset$ by the definition of Occ^* . When $X = YZ$, $Occ^*(X) = true$ if and only if

$$I(Y.u) \cap I_{\text{rev}}((Z.u)^R) \neq \emptyset,$$

by Lemma 6. The proof is complete. ■

Lemma 8 *For a simple collage system, we can build in $O(k^2\|\mathcal{D}\|+km)$ time using $O(k^2\|\mathcal{D}\|)$ space the table which stores the values $Output_1(t)$ for the tokens t in \mathcal{D} .*

Proof. When $X = a$, $Output_1(X) = true$ if and only if $D(P, a) \leq k$. When $X = YZ$,

$$Output_1(X) = Output_1(Y) \text{ OR } Output_1(Z) \text{ OR } Occ^*(X),$$

where OR denotes the logical OR. The proof is complete. ■

Now we turn into the computation of $Output_2$. The computation can be done with the same idea used for that of $Output_1$. By the definition of $Output_2$, it is not hard to see that $Output_2(S, t) = true$ if and only if

$$S \cap I_{\text{rev}}((t.u)^R) \cap \{\langle i, j \rangle \in Q \mid i \neq m\} \neq \emptyset.$$

Thus we can prove the next lemma.

Lemma 9 *For a simple collage system, the function $Output_2$ can be built in $O(k^2\|\mathcal{D}\|+km)$ time using $O(k^2\|\mathcal{D}\|)$ space, so that it returns its value in $O(k^2)$ time.*

Since $Output(S, t) = Output_1(t) \text{ OR } Output_2(S, t)$, we have the next result.

Lemma 10 *For a simple collage system, the function $Output$ can be built in $O(k^2\|\mathcal{D}\|+km)$ time using $O(k^2\|\mathcal{D}\|)$ space, so that it returns the value in $O(k^2)$ time.*

Theorem 2 follows from Lemmas 5 and 10.

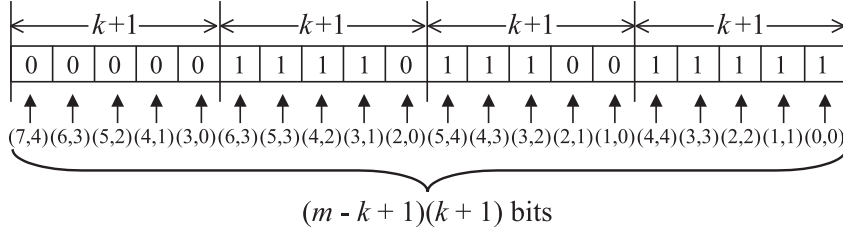


Figure 4: Encoding of the example NFA, where $m=7$ and $k=4$.

5.3 Bit-parallel implementation

In our implementation, we do not encode all of the automaton states in Q . We only encode the states in Q_{core} into a machine word. We thus substitute Q_{core} for Q . The discussion in Sections 5.1 and 5.2 is still true except for the computation of the table \bar{M} mentioned in the proof of Lemma 4. For a correct computation of the table \bar{M} , we need some modification based on the following fact.

Proposition 5 *Let J be a subset of $\{0, \dots, m\}$, and let $S = \{\langle i, j \rangle \in Q \mid i \in J\}$. Let $S_{\text{core}} = S \cap Q_{\text{core}}$. Then,*

$$Q_{\text{core}} \cap (S \oplus \langle 1, 0 \rangle) = Q_{\text{core}} \cap (S_{\text{core}} \oplus \langle 1, 0 \rangle \cup S_{\text{core}} \oplus \langle 1, 1 \rangle).$$

Now, we explain the details of our bit-parallel implementation. We use the $(k(i - j) + i)$ th bit of a machine word to indicate whether or not the state $\langle i, j \rangle \in Q_{\text{core}}$ is active (see Fig. 4). The union and the intersection operations on the subsets of Q_{core} can be evidently computed in $O(1)$ time by using the bitwise OR and AND, respectively. Next we discuss the computation of the operation \oplus . For a computation of $S \oplus \langle \ell, d \rangle$, it suffices to move all of the points $\langle i, j \rangle \in S$ to the right by $\ell - d$ (to the left by $d - \ell$, if $\ell < d$) and then move them down by d along the diagonal. The translation seems to be done using bit-shift by $(k + 1)(\ell - d) + d$. However, this is not enough because the bits in the result which correspond to the points $\langle i, j \rangle$ with $0 \leq j < d$ should be \emptyset . To make these bits 0, we use mask bits which depend on d and can be in advance built in $O(k)$ time and space. Therefore we can execute the operation \oplus in $O(1)$ time.

6 Conclusion

We have presented the first algorithm using bit-parallelism to the problem of approximate string matching on compressed text. The algorithm is intended to solve the existence problem. Assuming that our problem fits in a single machine word,

the algorithm runs in $O(k^2(\|\mathcal{D}\| + |\mathcal{S}|) + km)$ time using $O(k^2\|\mathcal{D}\|)$ space for a simple collage system. For the LZW or the LZ78 compression, the time and the space complexities are $O(k^2n + km)$ and $O(k^2n)$, respectively, where n is the compressed text length. This is a desirable property compared with the $O(mkn)$ time and space complexities of the algorithm due to Kärkkäinen et al. [4], although ours cannot cope with the all occurrence problem.

We have implemented our algorithm in a rather straightforward way. Unfortunately, the result of preliminary experiments on LZW compressed text shows that our current implementation of the algorithm is not better than a decompression followed by a search with the classical dynamic programming except for a large m and a small k . If we specialize the implementation to the LZW compression, the performance of our algorithm will be improved.

A practical fast algorithm for approximate string matching uses some filtering technique [7]. Although we have not yet devised a filtering technique suitable for compressed text searching, it will improve the performance of our algorithm in practice.

Acknowledgement

We would like to thank Professor Apostolico for fruitful discussions on this work.

References

- [1] A. Amir and G. Benson. Efficient two-dimensional compressed matching. In *Proc. Data Compression Conference*, page 279, 1992.
- [2] R. Baeza-Yates and G. Navarro. Faster approximate string matching. *Algorithmica*, 23:127–158, 1999.
- [3] J. E. Hopcroft and J. D. Ullman. *Introduction to Automata Theory, Languages, and Computation*. Addison-Wesley, 1979.
- [4] J. Kärkkäinen, G. Navarro, and E. Ukkonen. Approximate string matching over Ziv-Lempel compressed text. In *Proc. 11th Ann. Symp. on Combinatorial Pattern Matching*. Springer-Verlag, 2000. to appear.
- [5] T. Kida, Y. Shibata, M. Takeda, A. Shinohara, and S. Arikawa. A unifying framework for compressed pattern matching. In *Proc. 6th International Symp.*

- on String Processing and Information Retrieval*, pages 89–96. IEEE Computer Society, 1999.
- [6] N. J. Larsson and A. Moffat. Offline dictionary-based compression. In *Proc. Data Compression Conference '99*, pages 296–305. IEEE Computer Society, 1999.
- [7] G. Navarro. A guided tour to approximate string matching. Technical Report TR/DCC-99-5, Dept. of Computer Science, Univ. of Chile, 1999.
- [8] G. Navarro and M. Raffinot. A general practical approach to pattern matching over Ziv-Lempel compressed text. In *Proc. 10th Ann. Symp. on Combinatorial Pattern Matching*, pages 14–36. Springer-Verlag, 1999.
- [9] P. Sellers. The theory and computation of evolutionary distances: pattern recognition. *J. of Algorithms*, 1:359–373, 1980.
- [10] T. A. Welch. A technique for high performance data compression. *IEEE Comput.*, 17:8–19, June 1984.
- [11] S. Wu and U. Manber. Fast text searching allowing errors. *Comm. ACM*, 35(10):83–91, October 1992.
- [12] J. Ziv and A. Lempel. A universal algorithm for sequential data compression. *IEEE Trans. on Inform. Theory*, IT-23(3):337–349, May 1977.
- [13] J. Ziv and A. Lempel. Compression of individual sequences via variable-rate coding. *IEEE Trans. on Inform. Theory*, 24(5):530–536, Sep 1978.