



Scaled and permuted string matching

Ayelet Butman^a, Revital Eres^{b,1}, Gad M. Landau^{b,c,*,2}

^a *Holon Academic Institute of Technology, Israel*

^b *Department of Computer Science, Haifa University, Haifa 31905, Israel*

^c *Department of Computer and Information Science, Polytechnic University, Six MetroTech Center, Brooklyn, NY 11201-3840, USA*

Received 24 June 2004; received in revised form 6 September 2004

Available online 12 October 2004

Communicated by L.A. Hemaspaandra

Abstract

The goal of *scaled permuted string matching* is to find all occurrences of a pattern in a text, in all possible scales and permutations. Given a text of length n and a pattern of length m we present an $O(n)$ algorithm.

© 2004 Elsevier B.V. All rights reserved.

Keywords: Algorithms; Approximate string matching; Permutations; Scalings

1. Introduction

The well-known string matching problem that appears in all algorithm textbooks has as its input a text T of length n and a pattern P of length m over a given alphabet Σ . The output is all text locations where there is an exact match of the pattern. This problem has received much attention, and many algorithms have been developed to solve it (e.g., [6,13,14]). A detailed mod-

ern view of stringology can be found in a number of published books [5,7,11].

Most recent work has dealt with inexact matches. Many types of differences between the patterns were defined, for example, errors (Hamming distance, LCS [12], Edit distance [15]), rotations [1,3,9,10], scaling [2,4], or permutation. Most of the theoretical work has dealt with one type of difference at a time. This paper is one of the first attempts to deal with two types of differences together—scaling and permutation.

Definition 1 (*Scaled permuted string matching*).

Input: A pattern $P = p_1 \cdots p_m$ and a text $T = t_1 \cdots t_n$ both over alphabet Σ .

Output: All positions in T where an occurrence of a permuted copy of the pattern P , scaled to k starts

* Corresponding author.

E-mail addresses: butmosh@zahav.net.il (A. Butman), revitale@cs.haifa.ac.il (R. Eres), landau@poly.edu (G.M. Landau).

¹ Partially supported by the Israel Science Foundation grant 282/01.

² Partially supported by NSF grant CCR-0104307, by the Israel Science Foundation grant 282/01, and by IBM Faculty Award.

($k = 1, \dots, \lfloor n/m \rfloor$). The pattern is first permuted and then scaled.

Example. The string *bbbbaabbaaccaacc* is a scaled (to 2) permutation of *baabbacc*.

The scaled (only) string matching problem is a well studied problem. The algorithm presented in [4], which follows the method described in [8], achieves a linear running time for the scaled string matching problem. In [2] the case where the scaling of the pattern is by real numbers was considered, and a linear time algorithm was introduced.

An algorithm for the permuted string matching problem over run-length encoded strings is described in Section 2. In Section 3 we present the main result of this paper, an algorithm that solves the scaled permuted string matching problem in $O(n)$ time and space. Open problems are given in Section 4.

2. Permuted string matching over run-length encoded text

The permuted string matching problem over un-compressed text is simply solved. A sliding window of size $|P|$ can be moved over T to count, for each location of T , the order of statistics of the characters. Obviously, this can be done in $O(n)$ time.

The run-length of a string S is a popular encoding method. According to this encoding S can be described as a sequence of ordered pairs (σ, i) , often denoted by the *symbol* σ^i , each consisting of an alphabet *character* σ and an integer i . Each pair corresponds to a run in S , consisting of i consecutive occurrences of σ .

Let T' be the run-length compressed version of T where $T' = \sigma_1^{r_1} \dots \sigma_{|T'|}^{r_{|T'|}}$. Similarly, P' is the run-length compressed pattern. The pattern can be permuted, and therefore, in each location of the text we check if the order of statistics of the characters is equal to that of the pattern. As a result, a better compression can be achieved. Symbols with the same character are compressed. For example, let $P = aabbbaaccaab$, its run-length compressed version is $P' = a^2b^3a^1c^2a^2b^1$ and a permuted run-length compressed version is $P'' = a^5b^4c^2$. The technique we use is similar to the sliding window technique: a window

is shifted on T' from left to right in order to locate all the matches. The window is a substring of T' that represents a candidate for a match. Unlike the simple algorithm, this time the window size is not fixed.

We will define a *valid* window as a substring of T' that fulfills the following two properties:

sufficient The number of times each character appears in the window is at least the number of times it appears in the pattern.

minimal Removing the rightmost or the leftmost symbol of the window violates the *sufficient* property.

Note that:

- (a) The *valid* window property does not ensure a match.
- (b) If a permutation of the pattern occurs in a *valid* window of T' , $\sigma_i^{r_i} \dots \sigma_j^{r_j}$, then only the characters σ_i and σ_j can appear more times in this window than they appear in P'' .
- (c) If $\sigma_i = \sigma_j$ then the pattern may occur more than once in the *valid* window. Also, if $\sigma_{i-1} = \sigma_j$ ($\sigma_i = \sigma_{j+1}$) the pattern may occur more than once in $\sigma_{i-1}^{r_{i-1}} \dots \sigma_j^{r_j}$ ($\sigma_i^{r_i} \dots \sigma_{j+1}^{r_{j+1}}$).
- (d) A permuted pattern occurs in the text only in a *valid* window (including the symbols on the left and right of the window).

The algorithm scans the text, locates all *valid* windows and finds the ones in which a permuted copy of the pattern occurs. During the scan of the text, given a *valid* window, it is trivial to check if it contains a match. Hence, we will describe only how to locate all *valid* windows.

Note that given a text $T' = \sigma_1^{r_1} \dots \sigma_{|T'|}^{r_{|T'|}}$:

- (a) At most one *valid* window may start on each $\sigma_i^{r_i}$.
- (b) A *valid* window does not contain another *valid* window.

The *valid* windows are found by scanning the text from left to right, using two pointers, *left* and *right*. To discover each *valid* window, the *right* pointer moves first to find a *sufficient* window and then the *left* pointer moves to find the *valid* window within the *sufficient* window. Each move of the *right* pointer increases the

size of the window. The right pointer moves as long as deleting the leftmost symbol of the window violates the *sufficient* property of the window. When this symbol can finally be removed, the *right* pointer stops and the *left* pointer starts moving. Each move of the *left* pointer decreases the size of the window. The pointer moves as long as deleting the leftmost symbol of the window does not violate the *sufficient* property of the window. At this point, a new *valid* window has been found.

Example. Let $P'' = a^2b^3c^2d^2$ and $T' = c^3a^2c^2a^3d^2b^3c^1$ then $c^3a^2c^2a^3d^2b^3$ is the first *sufficient* window, and $c^2a^3d^2b^3$ is the first *valid* window (but not a match).

Claim 1. *The algorithm finds all (and only) valid windows.*

Proof. The algorithm reports only *valid* windows. We will prove by contradiction that the algorithm finds all the *valid* windows. Denote by Π_1 and Π_2 two consecutive *valid* windows that are discovered by the algorithm, and by i_{left_1} , i_{right_1} , i_{left_2} and i_{right_2} the left and right pointers of those windows, respectively. Assume that there exists a *valid* window Π_3 (with left and right pointers i_{left_3} and i_{right_3} , respectively) between Π_1 and Π_2 ($i_{\text{right}_1} < i_{\text{right}_3} < i_{\text{right}_2}$) that the algorithm does not discover. By the *minimal* propriety we get that $i_{\text{left}_1} < i_{\text{left}_3}$. After reporting Π_1 the algorithm looks for the next *valid* window. During the scanning of the right pointer the algorithm passes i_{right_3} and does not stop, which means that the window $i_{\text{left}_1} + 1 \cdots i_{\text{right}_3}$ does not satisfy the sufficient property. Since, $i_{\text{left}_1} + 1 \leq i_{\text{left}_3}$ we conclude that the window Π_3 does not satisfy the sufficient property as well, and hence, it is not *valid*. \square

Time complexity. We assume that $|\Sigma|$ is $O(|P''|)$, hence, the time complexity of the algorithm is $O(|P''| + |T'|)$. In case the input pattern is not given in a permuted run-length compressed format, an $O(|P|)$ time preprocessing step is added.

3. A linear time algorithm for the scaled permuted string matching problem

The algorithm is composed of two stages:

- (1) Preprocessing the text T' . Computing compact copies of the text for each possible scale $1 \leq s \leq n/m$.
- (2) Applying the permuted string matching over the run-length encoded text algorithm (Section 2) on the copies of the text.

Observation 1. If a permutation of P scaled to s occurs in $\sigma_i^{j_i} \cdots \sigma_k^{j_k}$ then j_{i+1}, \dots, j_{k-1} are multiples of s , and $j_i, j_k \geq s$.

Following the above observation, we compute for each scale s a compact text T'_s in the following two steps:

Step 1: Locate all the regions in T' where the symbols appear with multiples of s . Add the symbol \$ as a separator between the regions.

Step 2: Expand these regions to include the symbols on their boundaries. In order to simplify the computation of stage 2, a symbol $t_j^{r_j}$ of T' is replaced in T'_s by $t_j^{\lfloor r_j/s \rfloor}$.

Step 1. Locating the regions— T' is scanned from left to right. Consider a symbol $t_i^{r_i}$. A new symbol $t_i^{r_i/s}$ is added to T'_s if r_i is a multiple of s . The following code describes this idea.

Step 1—The parallel construction of the new text

For every symbol in T' do:

{ let a^r be the current symbol being examined }

$s = 1$

Repeat Until $s > \sqrt{r}$

If $(r \bmod s = 0)$ Then

Add $a^{r/s}$ to T'_s

{ skip the next line if $s = \sqrt{r}$ }

Add a^s to $T'_{r/s}$

$s = s + 1$

Note that the efficiency of this procedure depends on the method that finds all the divisors of an integer. In the above example we used a naive method. A new symbol that is added at the end of T'_s may continue a region or start a new one. In the second case we add a separator (\$) between the regions.

Step 2. Expansion of the regions—The last refinement is done by scanning each T'_s text from left to

right and expanding all the regions we generated in step 1. In the next procedure we deal with symbols that appear on the left side of a \$ separator in T'_s . The opposite case is treated in the same way.

Step 2

For every \$ separator in T'_s do:

{ let $t_i^{r_i/s}$ be the symbol appearing on the left side of the current \$ separator on T'_s , and let $t_{i+1}^{r_{i+1}}$ be the adjacent symbol to $t_i^{r_i}$ on T' }

If $(r_{i+1} > s)$ then

Add $t_{i+1}^{\lfloor (r_{i+1})/s \rfloor}$ to T'_s between $t_i^{r_i/s}$ and the \$ separator

Example. Let $T' = a^6b^2c^4a^3d^5b^9d^2c^8b^4a^7$, the new text after applying step 2 is:

$$T'_1 = \$a^6b^2c^4a^3d^5b^9d^2c^8b^4a^7$,$$

$$T'_2 = \$a^3b^1c^2a^1$b^4d^1c^4b^2a^3$,$$

$$T'_3 = \a^2c^1a^1d^1d^1b^3,$$

$$T'_4 = \c^1c^2b^1a^1$,$$

$$T'_5 = d^1b^1,$$

$$T'_6 = a^1,$$

$$T'_7 = a^1,$$

$$T'_8 = c^1,$$

$$T'_9 = b^1.$$

Stage 2 runs the permuted string matching over a run-length encoded text algorithm (Section 2) on all the new compact texts.

Time complexity. The input to our problem is a compressed text $T' = t_1^{r_1} \dots t_k^{r_k}$, whose original length is n , and a pattern P'' of length $|P''|$ (or a pattern P of length m). Both the pattern and the text are over alphabet Σ . The following claim shows that the total length of all compact new texts is linear.

Claim 2. *The total length of all the new texts T'_s ($1 \leq s \leq n/m$) is $O(n)$.*

Proof. In step 1, we consider each symbol $t_i^{r_i}$ in T' , and the number of new symbols that we produce from $t_i^{r_i}$ is bounded by $2\sqrt{r_i}$. In addition we may add a \$

separator to each new symbol. In step 2, two new symbols may be added to each \$ separator. Hence the total length of all new texts is: $8 \cdot \sum_{i=1}^k \sqrt{r_i} = O(n)$. \square

The running time of both stage 1 and stage 2 is bounded by the length of the new texts, hence the total time complexity is $O(n)$.

4. Open problems

The algorithm described in this paper is the first to deal with scaling and permutation. We considered the case in which the pattern is first permuted and then scaled. The first challenge is to design an $o(nm)$ algorithm for the case in which the pattern is first scaled and then permuted. We also dealt with integer scales. The second challenge is to deal with scales that are real numbers.

Acknowledgements

The authors are grateful to the referees for their helpful comments.

References

- [1] A. Amir, A. Butman, M. Crochemore, G.M. Landau, M. Schaps, Two-dimensional pattern matching with rotations, *Theoret. Comput. Sci.* 314 (1–2) (2004) 173–187.
- [2] A. Amir, A. Butman, M. Lewenstein, Real scaled matching, *Inform. Process. Lett.* 70 (4) (1999) 185–190.
- [3] A. Amir, O. Kapah, D. Tsur, Faster two dimensional pattern matching with rotations, in: *Proc. 15th Annual Symposium on Combinatorial Pattern Matching (CPM 2004)*, in: *Lecture Notes in Comput. Sci.*, vol. 3109, Springer, Berlin, 2004.
- [4] A. Amir, G.M. Landau, U. Vishkin, Efficient pattern matching with scaling, *J. Algorithms* 13 (1) (1992) 2–32.
- [5] A. Apostolico, Z. Galil (Eds.), *Pattern Matching Algorithms*, Oxford University Press, 1997.
- [6] R.S. Boyer, J.S. Moore, A fast string searching algorithm, *Comm. ACM* 20 (1977) 762–772.
- [7] M. Crochemore, W. Rytter, *Text Algorithms*, Oxford University Press, 1994.
- [8] T. Eilam-Tsoreff, U. Vishkin, Matching patterns in strings subject to multilinear transformations, *Theoret. Comput. Sci.* 60 (3) (1988) 231–254.
- [9] K. Fredriksson, G. Navarro, E. Ukkonen, Optimal exact and fast approximate two dimensional pattern matching allowing rotations, in: *Proceedings of the 13th Annual Symposium*

- on Combinatorial Pattern Matching (CPM 2002), in: *Lecture Notes in Comput. Sci.*, vol. 2373, Springer, Berlin, 2002, pp. 235–248.
- [10] K. Fredriksson, E. Ukkonen, A rotation invariant filter for two-dimensional string matching, in: *Proc. 9th Annual Symposium on Combinatorial Pattern Matching (CPM 1998)*, in: *Lecture Notes in Comput. Sci.*, vol. 1448, Springer, Berlin, 1998, pp. 118–125.
- [11] D. Gusfield, *Algorithms on Strings, Trees, and Sequences: Computer Science and Computational Biology*, Cambridge University Press, 1997.
- [12] D.S. Hirshberg, Algorithms for the longest common subsequence problem, *J. ACM.* 24 (4) (1977) 664–675.
- [13] R. Karp, M.O. Rabin, Efficient randomized pattern-matching algorithms, *IBM J. Res. Dev.* (1987) 249–260.
- [14] D.E. Knuth, J.H. Morris, V.R. Pratt, Fast pattern matching in strings, *SIAM J. Comput.* 6 (1977) 323–350.
- [15] V.I. Levenshtein, Binary codes capable of correcting, deletions, insertions and reversals, *Soviet Phys. Dokl.* 10 (1966) 707–710.