

Kybernetika

VOLUME 38 (2002), NUMBER 1

The Journal of the Czech Society for
Cybernetics and Information Sciences

Published by:

Institute of Information Theory
and Automation of the Academy
of Sciences of the Czech Republic

Editor-in-Chief:

Milan Mareš

Managing Editors:

Karel Sladký

Editorial Board:

Jiří Anděl, Marie Demlová, Petr Hájek,
Jan Hlavička, Martin Janžura, Jan Ježek,
Radim Jiroušek, Ivan Kramosil,
František Matúš, Jiří Outrata, Jan Štecha,
Olga Štěpánková, Igor Vajda, Pavel Zítek,
Pavel Žampa

Editorial Office:

Pod Vodárenskou věží 4, 182 08 Praha 8

Kybernetika is a bi-monthly international journal dedicated for rapid publication of high-quality, peer-reviewed research articles in fields covered by its title.

Kybernetika traditionally publishes research results in the fields of Control Sciences, Information Sciences, System Sciences, Statistical Decision Making, Applied Probability Theory, Random Processes, Fuzziness and Uncertainty Theories, Operations Research and Theoretical Computer Science, as well as in the topics closely related to the above fields.

The Journal has been monitored in the Science Citation Index since 1977 and it is abstracted/indexed in databases of Mathematical Reviews, Current Mathematical Publications, Current Contents ISI Engineering and Computing Technology.

Kybernetika. Volume 38 (2002)

ISSN 0023-5954, MK ČR E 4902.

Published bi-monthly by the Institute of Information Theory and Automation of the Academy of Sciences of the Czech Republic, Pod Vodárenskou věží 4, 182 08 Praha 8. — Address of the Editor: P. O. Box 18, 182 08 Prague 8, e-mail: kybernetika@utia.cas.cz. — Printed by PV Press, Pod vrstevnicí 5, 140 00 Prague 4. — Orders and subscriptions should be placed with: MYRIS TRADE Ltd., P. O. Box 2, V Štíhlách 1311, 142 01 Prague 4, Czech Republic, e-mail: myris@myris.cz. — Sole agent for all “western” countries: Kubon & Sagner, P. O. Box 34 01 08, D-8 000 München 34, F.R.G.

Published in February 2002.

© Institute of Information Theory and Automation of the Academy of Sciences of the Czech Republic, Prague 2002.

A NEW PRACTICAL LINEAR SPACE ALGORITHM FOR THE LONGEST COMMON SUBSEQUENCE PROBLEM*

HEIKO GOEMAN AND MICHAEL CLAUSEN

This paper deals with a new practical method for solving the longest common subsequence (LCS) problem. Given two strings of lengths m and n , $n \geq m$, on an alphabet of size s , we first present an algorithm which determines the length p of an LCS in $O(ns + \min\{mp, p(n-p)\})$ time and $O(ns)$ space. This result has been achieved before [29, 30], but our algorithm is significantly faster than previous methods. We also provide a second algorithm which generates an LCS in $O(ns + \min\{mp, m \log m + p(n-p)\})$ time while preserving the linear space bound, thus solving the problem posed in [29, 30]. Experimental results confirm the efficiency of our method.

1. INTRODUCTION

Let $x = x_1 \dots x_m$ and $y = y_1 \dots y_n$, $n \geq m$, be two strings over an alphabet $\Sigma = \{\sigma_1, \dots, \sigma_s\}$ of size s . A *subsequence* of x is a sequence of symbols obtained by deleting zero or more characters from x . The *Longest Common Subsequence (LCS) Problem* is to find a common subsequence of x and y which is of greatest possible length.

It will be convenient to describe the problem in another way. An ordered pair (k, ℓ) , $1 \leq k \leq m$, $1 \leq \ell \leq n$, is called a *match* if $x_k = y_\ell$. The set M of all matches can be identified with a *matching matrix* of size $m \times n$ in which each match is marked with a dot. For example, if $x = abacbcba$ and $y = cbabbacac$, then M is as shown in Figure 1 (a). Define a partial order \ll on $\mathbb{N} \times \mathbb{N}$ by establishing $(k, \ell) \ll (k', \ell')$ iff both $k < k'$ and $\ell < \ell'$. A *chain* $C \subseteq M$ is a set of points which are pairwise comparable, i. e., for any two distinct $p_1, p_2 \in C$, either $p_1 \ll p_2$ or $p_1 \gg p_2$, where $p_1 \gg p_2$ means $p_2 \ll p_1$. Then the LCS problem can be viewed as finding a chain of maximal cardinality in M . One such chain is indicated as a path in Figure 1 (b).

Finding an LCS is closely related with the computation of string edit distances [21, 24, 34, 36] and shortest common supersequences [14]. It was first used by biologists to study amino acids [9, 10, 27, 31]. Other applications are in data compression [1, 14, 23] and pattern recognition [13, 22].

*Research supported by Deutsche Forschungsgemeinschaft, Grant CL 64/3-1.

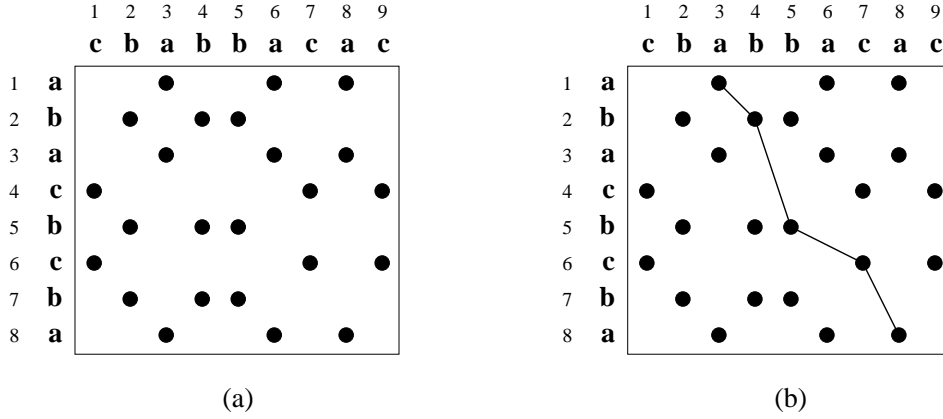


Fig. 1. (a) matching matrix, (b) path representing an LCS.

The LCS problem can be solved in $O(mn)$ time by a dynamic programming approach [32, 35], while the asymptotically fastest general solution uses the “four russians” trick and takes $O(nm/\log n)$ time [24]. A lot of other algorithms have also been developed which are sensitive to other problem parameters, e. g., the length p of an LCS. They usually perform much better than the latter algorithms, although they all have a worst case time complexity at least of $\Omega(mn)$. To give an example, Hunt and Szymanski [19] have presented an $O((r + n) \log n)$ algorithm, where $r := |M|$. Thus their approach is fast when r is small, e. g., $r = O(n)$, but its worst-case time complexity is $O(n^2 \log n)$. Later, this has been improved to $O(mn)$ [2]. There are also several routines [25, 26, 33, 37] which run in $O(n(n + 1 - p))$ or $O(n(m + 1 - p))$ time, and thus are efficient when an LCS is expected to be long. Other algorithms have running times $O(n(p + 1))$ or $O(m(p + 1))$ and should be used for short LCS [3, 4, 17, 18]. However, it might be very difficult to *a priori* select a good strategy because in general the length p cannot be easily estimated. Also, when having a small alphabet, we can expect p to be of intermediate size, e. g., for $s = 4$, the average length of an LCS is bounded between $0.54 \cdot m \leq p \leq 0.71 \cdot m$ [7, 8, 11, 28, 32]. Then none of the above methods performs well. Therefore recent research has been concentrated on more flexible algorithms which are efficient for short, intermediate, and long LCS, such as the method proposed by Chin/Poon [6]. Another approach from Rick [29, 30] with running time $O(ns + \min\{mp, p(n - p)\})$ has been widely accepted as the fastest algorithm for the general LCS problem.

In this paper, we shall develop a new algorithm which is based on a kind of dualization of Rick’s method. A detailed description of the theoretical background will be given in Sections 2 and 3. We do not improve the $O(ns + \min\{mp, p(n - p)\})$ time bound, but two important advantages are obtained. First, the number of matches processed while computing the length of an LCS is significantly decreased, resulting in a faster execution speed. The corresponding algorithm will be presented in Section 4. Second, when generating an LCS, we can achieve linear space through

a divide-and-conquer scheme similar to that of several other (but slower) algorithms [5, 16, 20]. This will be explained in Section 5. The methods mentioned before all need at least $\Omega(nm/\log n)$ space in their worst cases (see [28] for a survey), and most of them, including Rick's approach, cannot be combined with the divide-and-conquer technique. The open problem of a linear space implementation of Rick's algorithm [30] is hereby solved. Experimental results presented in Section 6 confirm the efficiency of our method.

2. A NEW APPROACH TO THE LCS PROBLEM

As already mentioned in the introduction, the LCS problem is equivalent to finding a chain of maximum cardinality in M . Dilworth's fundamental theorem [12] states that this cardinality equals the minimum number of disjoint *antichains* into which M can be decomposed (an antichain of M consists of matches which are pairwise incomparable). In our example, this number (called the *Sperner number* of M) equals five. A suitable decomposition is shown in Figure 2(f). To find such a minimum decomposition, we first split $[1 : m] \times [1 : n]$ into subsets denoted by T^i , L^i , B^i , and R^i , where

$$\begin{aligned} T^i &:= \{i\} \times [i : n + 1 - i] \\ L^i &:= [i + 1 : m + 1 - i] \times \{i\} \\ B^i &:= \{m + 1 - i\} \times [i + 1 : n + 1 - i] \\ R^i &:= [i + 1 : m - i] \times \{n + 1 - i\} \end{aligned}$$

and $1 \leq i \leq \lceil m/2 \rceil$ (see Figure 2(a) for an illustration). Additionally, let

$$T^{\leq i} := \bigcup_{j \leq i} T^j, \quad L^{\leq i} := \bigcup_{j \leq i} L^j, \quad B^{\leq i} := \bigcup_{j \leq i} B^j, \quad R^{\leq i} := \bigcup_{j \leq i} R^j.$$

Now for $i = 1, 2, \dots, \lceil m/2 \rceil$, we construct four sets of antichains $A^{T,i}$, $A^{L,i}$, $A^{B,i}$, and $A^{R,i}$ which decompose (a suitable subset of) $T^{\leq i}$, $L^{\leq i}$, $B^{\leq i}$, and $R^{\leq i}$, respectively. The decompositions are generated by updating the previous sets, using the matches found in T^i , L^i , B^i , and R^i (details are given below). We use $A_u^{T,i}$ to denote an antichain in $A^{T,i}$, where u is an index between 1 and the size $e^{T,i} := |A^{T,i}|$ of $A^{T,i}$. Therefore $e^{T,i}$ is also called the *end index* of $A^{T,i}$. For $A^{L,i}$, $A^{B,i}$, and $A^{R,i}$, we introduce analogous notations. Furthermore, there are two *start indices* $s^{TL,i}$ and $s^{BR,i}$. The first one is used to split both $A^{T,i}$ and $A^{L,i}$ into two parts. One part contains all antichains with indices less than $s^{TL,i}$, and the other part consists of the rest. Only the latter part will be used for the updating process, whereas the former one will be copied to $A^{T,i+1}$ resp. $A^{L,i+1}$ without change. $s^{BR,i}$ similarly splits $A^{B,i}$ and $A^{R,i}$.

Figure 2(b), (c), (d), and (e) give a preview of the construction in the sample matching matrix after step $i = 1, 2, 3$, and 4, respectively. The centered grey box represents the remaining part of M which has not been processed so far. By our construction, with each step, it shrinks by two rows and columns.

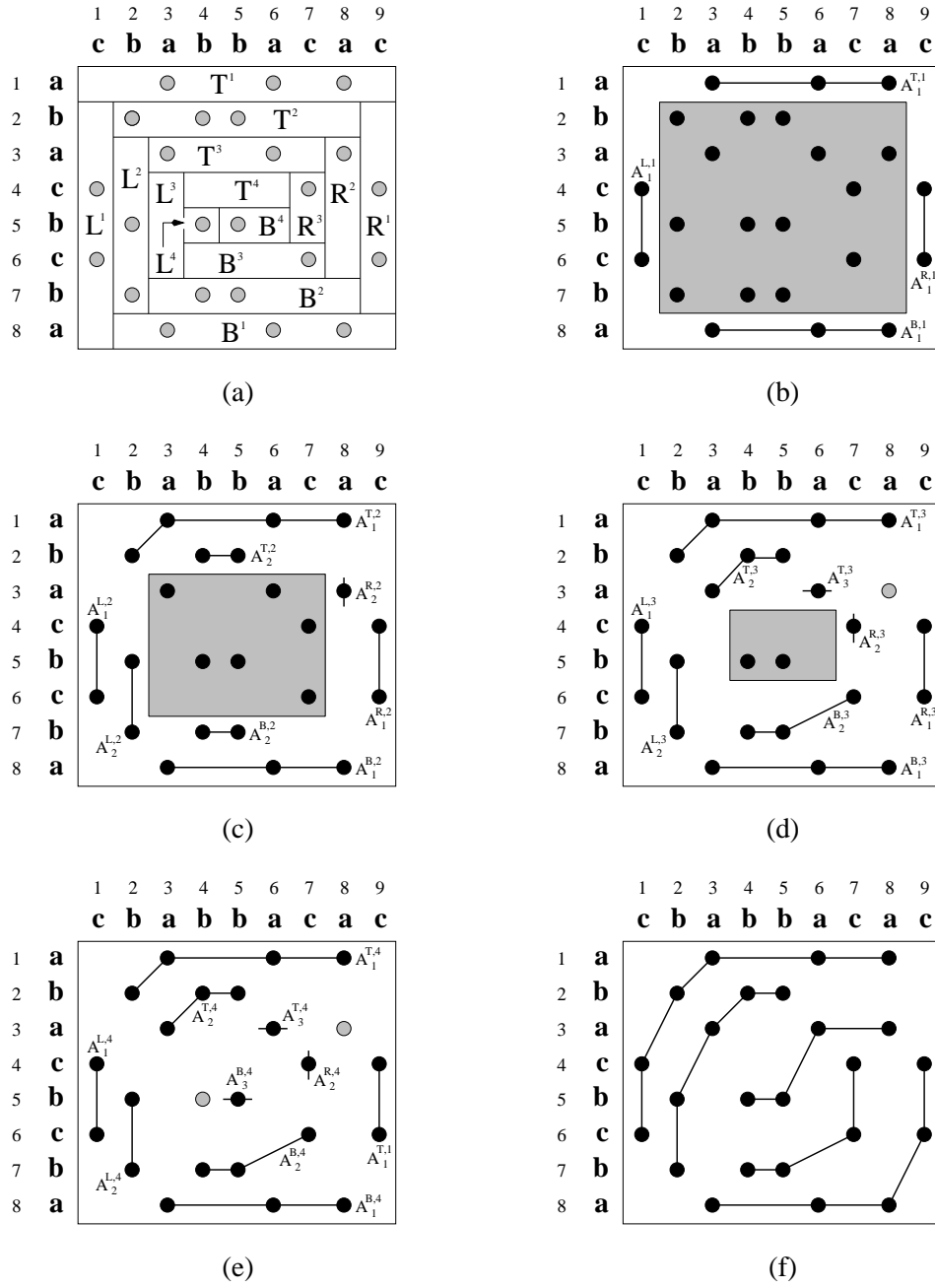


Fig. 2. (a) splitting of M , (b)–(e) construction of antichains, (f) final decomposition.

We need the following terminology for the description of the construction process. For two antichains $C, D \subseteq M$ the set

$$IP(C, D) := \{p_1 \in C \mid \forall p_2 \in D: \neg(p_1 \ll p_2 \vee p_1 \gg p_2)\}$$

is called the *incomparable part* of C relative to D . Clearly, $IP(C, D) \cup D$ is the greatest antichain above D contained in $C \cup D$. We say C is *incomparable* to D if $IP(C, D) = C$. Furthermore, a single match $p_1 \in M$ is *incomparable* to D if $IP(\{p_1\}, D) = \{p_1\}$.

We are now prepared to discuss the generation of the antichains in more detail. Initially, there are no antichains, i. e., we have $A^{T,0} = A^{L,0} = A^{B,0} = A^{R,0} = \emptyset$ by initializing each start and end index to 1 and 0, respectively. Then, for each step $i = 1, \dots, \lceil m/2 \rceil$, we start with T^i to determine $A^{T,i}$ from $A^{T,i-1}$. Let $s := s^{TL,i-1}$ and $e := e^{T,i-1}$. The first $s-1$ antichains remain unchanged and are simply copied from $A^{T,i-1}$ to $A^{T,i}$. Now define $A_s^{T,i}$ as $A_s^{T,i-1} \cup IP(T^i \cap M, A_s^{T,i-1})$. For example, when processing T^2 in Figure 2 (b), $IP(T^2 \cap M, A_1^{T,1}) = \{(2, 2)\}$, and thus the match $(2, 2)$ combined with $A_1^{T,1}$ makes up $A_1^{T,2}$ as shown in Figure 2 (c). Next, setting $u = s+1, \dots, e$, the antichain $A_u^{T,i-1}$ is handled in the same way to set up $A_u^{T,i}$, but only those matches in T^i not belonging to $A_s^{T,i}, \dots, A_{u-1}^{T,i}$ are considered. Finally, we establish $s^{TL,i} := s$ and, if there are no matches left, $e^{T,i} := e$. Otherwise, we set $e^{T,i}$ to $e+1$ and collect all remaining matches in a new antichain $A_{e+1}^{T,i}$. Also, if $A^{R,i-1} \neq \emptyset$, we check whether its last antichain $A_{\tilde{e}}^{R,i-1}$, $\tilde{e} := e^{R,i-1}$, is incomparable to $A_{e+1}^{T,i}$. In this case we say $A_{\tilde{e}}^{R,i-1}$ is *inactivated* by $A_{e+1}^{T,i}$, and we remove $A_{\tilde{e}}^{R,i-1}$ from $A^{R,i}$ by setting $e^{R,i} := e^{R,i-1}$. Continuing our example with T^2 in Figure 2 (b), we see there are two matches $(2, 4)$ and $(2, 5)$ left after processing $A_1^{T,2}$. Therefore a new antichain $A_2^{T,2}$ is created, but $A_1^{R,1}$ remains unchanged because, for example, $(2, 4) \ll (4, 9)$. The final set $A^{T,2}$ is shown in Figure 2 (c) (the modifications to the other antichains are described below). Now let us consider the work involved with T^3 . The match $(3, 3)$ cannot be put into $A_1^{T,3}$, but into $A_2^{T,3}$, and the other match $(3, 6)$ makes up the new antichain $A_3^{T,3}$. This time $(3, 6)$ inactivates $(3, 8)$, and thus $A_2^{R,2}$ is removed. The result is illustrated in Figure 2 (d) (all matches located in deleted antichains are indicated by grey dots).

Having determined $A^{T,i}$, we continue with the necessary calculations for $A^{L,i}$ which are very similar. Again, the first $s-1$ antichains are copied. Then, by setting $u = s, \dots, e^{L,i-1}$, $A_u^{L,i}$ is defined as the union of $A_u^{L,i-1}$ and the incomparable part of L^i relative to $A_u^{L,i-1}$, where only those matches are considered which have not already been used. Remaining matches form a new antichain and, if they are incomparable to the last antichain in $A^{B,i-1}$, we decrease $e^{B,i}$ by one. The corresponding algorithm in Figure 3 (a) also introduces two additional sets D^{TR} and D^{BL} which contain all deleted matches. Details will be given in the next section.

Before processing $A^{B,i-1}$ and $A^{R,i-1}$ in an analogous way, we first check whether the first antichain in $A^{T,i}$ or $A^{L,i}$ is *TL-complete*, i. e., whether one of them contains a match (k, ℓ) such that $1 \leq k, \ell \leq i$. For example, in the configuration shown in Figure 2 (c), $A_1^{T,2}$ is TL-complete due to the match $(2, 2)$. As soon as $A_s^{T,i}$ is detected to be TL-complete, $s^{TL,i}$ is increased by one, thus the first antichains in both corresponding sets which are checked for additional matches remain unchanged

(a)

```

S := T^i ∩ M;      (* Determine A^{T,i} *)
For u := s^{TL,i-1} To e^{T,i-1} Do {
  A_u^{T,i} := A_u^{T,i-1} ∪ IP(S, A_u^{T,i-1});
  S := S \ IP(S, A_u^{T,i-1});
5 };
If S ≠ ∅ Then {
  e^{T,i} := e^{T,i-1} + 1; e := e^{T,i}; A_e^{T,i} := S;
  e^{R,i} := e^{R,i-1}; ē := e^{R,i};
  If s^{BR,i-1} ≤ e^{R,i-1} Then {
10   If IP(A_e^{R,i-1}, A_e^{T,i}) = A_e^{R,i-1} Then {
     D^{TR} := D^{TR} ∪ A_e^{R,i-1};
     e^{R,i} := ē - 1;
   };
  };
15 } Else { e^{T,i} := e^{T,i-1}; e^{R,i} := e^{R,i-1} };
For u := 1 To s^{TL,i-1} - 1 Do A_u^{T,i} := A_u^{T,i-1};

S := L^i ∩ M;      (* Determine A^{L,i} *)
For u := s^{TL,i-1} To e^{L,i-1} Do {
  A_u^{L,i} := A_u^{L,i-1} ∪ IP(S, A_u^{L,i-1});
20 S := S \ IP(S, A_u^{L,i-1});
};
If S ≠ ∅ Then {
  e^{L,i} := e^{L,i-1} + 1; e := e^{L,i}; A_e^{L,i} := S;
  e^{B,i} := e^{B,i-1}; ē := e^{B,i};
25 If s^{BR,i-1} ≤ e^{B,i-1} Then {
   If IP(A_e^{B,i-1}, A_e^{L,i}) = A_e^{B,i-1} Then {
     D^{BL} := D^{BL} ∪ A_e^{B,i-1};
     e^{B,i} := ē - 1;
   };
  };
30 } Else { e^{L,i} := e^{L,i-1}; e^{B,i} := e^{B,i-1} };
For u := 1 To s^{TL,i-1} - 1 Do A_u^{L,i} := A_u^{L,i-1};
33 s^{TL,i} := s^{TL,i-1};

```

(b)

```

S := B^i ∩ M;      (* Determine A^{B,i} *)
For u := s^{BR,i-1} To e^{B,i-1} Do {
  A_u^{B,i} := A_u^{B,i-1} ∪ IP(S, A_u^{B,i-1});
  S := S \ IP(S, A_u^{B,i-1});
};
If S ≠ ∅ Then {
  e^{B,i} := e^{B,i-1} + 1; e := e^{B,i}; A_e^{B,i} := S;
  If s^{TL,i} ≤ e^{L,i-1} Then {
    ē := e^{L,i};
    If IP(A_e^{L,i-1}, A_e^{B,i}) = A_e^{L,i-1} Then {
      D^{BL} := D^{BL} ∪ A_e^{L,i-1};
      e^{L,i} := ē - 1;
    };
  };
};
For u := 1 To s^{BR,i-1} - 1 Do A_u^{B,i} := A_u^{B,i-1};

S := R^i ∩ M;      (* Determine A^{R,i} *)
For u := s^{BR,i-1} To e^{R,i-1} Do {
  A_u^{R,i} := A_u^{R,i-1} ∪ IP(S, A_u^{R,i-1});
  S := S \ IP(S, A_u^{R,i-1});
};
If S ≠ ∅ Then {
  e^{R,i} := e^{R,i-1} + 1; e := e^{R,i}; A_e^{R,i} := S;
  If s^{TL,i} ≤ e^{T,i-1} Then {
    ē := e^{T,i};
    If IP(A_e^{T,i-1}, A_e^{R,i}) = A_e^{T,i-1} Then {
      D^{TR} := D^{TR} ∪ A_e^{T,i-1};
      e^{T,i} := ē - 1;
    };
  };
};
For u := 1 To s^{BR,i-1} - 1 Do A_u^{R,i} := A_u^{R,i-1};
s^{BR,i} := s^{BR,i-1};

```

Fig. 3. The algorithms for generating $A^{T,i}$ & $A^{L,i}$ (a), and $A^{B,i}$ & $A^{R,i}$ (b).

from now on. If there is no such antichain in $A^{L,i}$ (i.e. $s > e^{L,i}$), but $s^{BR,i-1} \leq e^{B,i}$, then we additionally test whether $A_s^{T,i}$ is incomparable to the last antichain in $A^{B,i-1}$ and, should this situation arise, delete this antichain from $A^{B,i}$ by decreasing $e^{B,i}$.

Now assume $A_s^{L,i}$ is TL-complete. Then, as shown in Figure 4 (a), we also increase $s^{TL,i}$, and similarly, if $s > e^{T,i}$ and $s^{BR,i-1} \leq e^{R,i}$, we decrease $e^{R,i}$ if $A_s^{L,i}$ inactivates the last antichain in $A^{R,i}$.

The remaining work in step i concerns with the analogous construction of $A^{B,i}$ and $A^{R,i}$. (The analogue of TL-completeness is called *BR-completeness*. An antichain is BR-complete if it contains a match (k, ℓ) with $m - i < k \leq m$ and $n - i < \ell \leq n$.) Details are available from the algorithms shown in Figure 3 (b) and Figure 4 (b).

The main program shown in Figure 5 is straightforward. Our next task is to elaborate the connection between the generated antichains and a minimal decomposition of M . This is done in the next section.

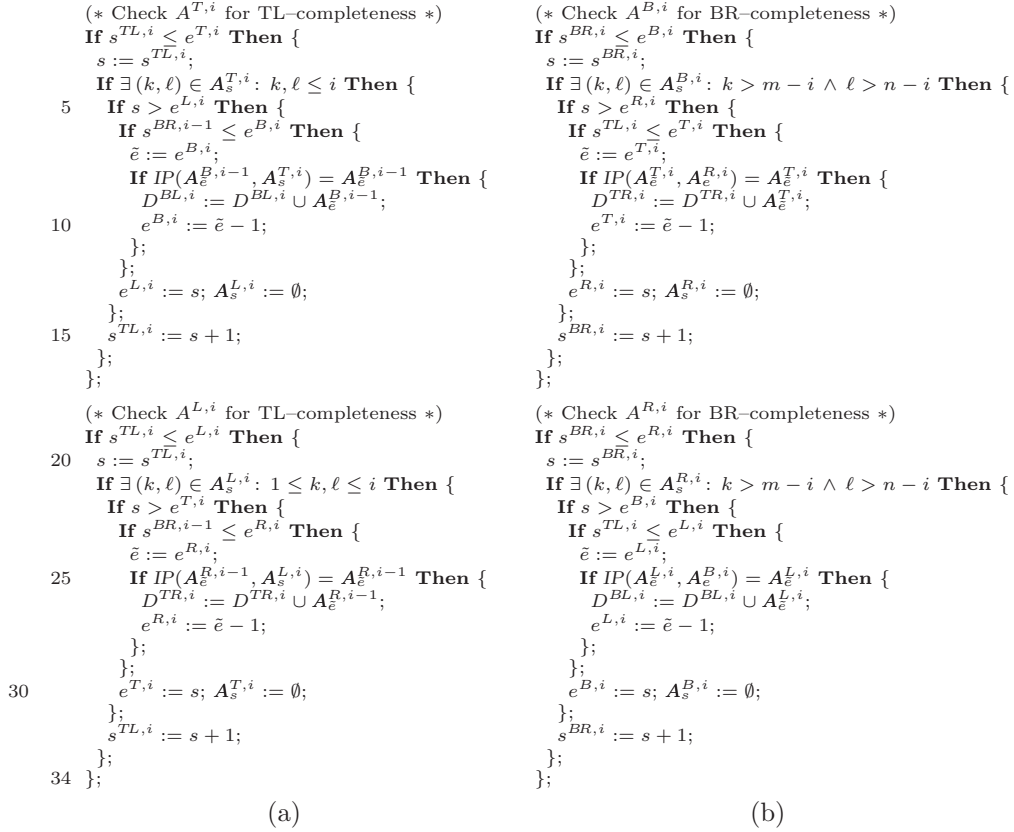


Fig. 4. The algorithms for handling complete antichains in $A^{T,i}$ & $A^{L,i}$ (a), and in $A^{B,i}$ & $A^{R,i}$ (b).

3. ANALYSIS OF THE CONSTRUCTION

In this section, we study how to combine the antichains into larger ones such that a minimal decomposition of M is obtained. We further establish some results which later help us to construct an LCS in linear space.

Let us assume m is odd, and let $i = \lceil m/2 \rceil$. For technical reasons, we then put $A_u^{B,i} := A_u^{B,i-1}$ and $A_u^{R,i} := A_u^{R,i-1}$ for all $1 \leq u \leq e^{B,i-1}$ and $1 \leq u \leq e^{R,i-1}$. We also set $s^{BR,i} := s^{BR,i-1}$, $e^{B,i} := e^{B,i-1}$, and $e^{R,i} := e^{R,i-1}$. Furthermore, for $0 \leq i \leq \lceil m/2 \rceil$, we define $A_u^{T,i} := \emptyset$, $A_u^{L,i} := \emptyset$, $A_u^{B,i} := \emptyset$, and $A_u^{R,i} := \emptyset$ for $u > e^{T,i}$, $u > e^{L,i}$, $u > e^{B,i}$, and $u > e^{R,i}$, respectively.

Lemma 3.1. Let $1 \leq i \leq \lceil m/2 \rceil$. Then the following holds:

- a) $\forall s^{TL,i-1} \leq u < v \leq e^{T,i} \forall p_1 \in A_v^{T,i} \exists p_2 \in A_u^{T,i} : p_1 \gg p_2$.
- b) $\forall s^{TL,i-1} \leq u < v \leq e^{L,i} \forall p_1 \in A_v^{L,i} \exists p_2 \in A_u^{L,i} : p_1 \gg p_2$.


```

 $s^{T,0} := 1; s^{L,0} := 1; s^{B,0} := 1; s^{R,0} := 1; \quad (* \text{ Initialization } *)$ 
 $e^{T,0} := 0; e^{L,0} := 0; e^{B,0} := 0; e^{R,0} := 0;$ 
For  $i := 0$  To  $\lceil m/2 \rceil$  Do  $D^{TL,i} := \emptyset;$ 
For  $i := 0$  To  $\lfloor m/2 \rfloor$  Do  $D^{BR,i} := \emptyset;$ 
5  $i := 1;$ 
While  $i \leq \lfloor m/2 \rfloor$  Do {  $(*)$  Main loop  $(*)$ 
    Determine  $A^{T,i}$  and  $A^{L,i}$ ;  $(*)$  see Figure 3 (a)  $(*)$ 
    Look for TL-complete antichains in  $A^{T,i}$  and  $A^{L,i}$ ;  $(*)$  see Figure 4 (a)  $(*)$ 
    Determine  $A^{B,i}$  and  $A^{R,i}$ ;  $(*)$  see Figure 3 (b)  $(*)$ 
10 Look for BR-complete antichains in  $A^{B,i}$  and  $A^{R,i}$ ;  $(*)$  see Figure 4 (b)  $(*)$ 
     $i := i + 1;$ 
};
If  $\text{Odd}(m)$  Then {
    Determine  $A^{T,\lceil m/2 \rceil}$  and  $A^{L,\lceil m/2 \rceil}$ ;  $(*)$  see Figure 3 (a)  $(*)$ 
15 Look for TL-complete antichains in  $A^{T,\lceil m/2 \rceil}$  and  $A^{L,\lceil m/2 \rceil}$ ;  $(*)$  see Figure 4 (a)  $(*)$ 
};

```

Fig. 5. The main program for decomposing M .

$$c) \forall s^{BR,i-1} \leq u < v \leq e^{B,i} \forall p_1 \in A_v^{B,i} \exists p_2 \in A_u^{B,i}: p_1 \ll p_2.$$

$$d) \forall s^{BR,i-1} \leq u < v \leq e^{R,i} \forall p_1 \in A_v^{R,i} \exists p_2 \in A_u^{R,i}: p_1 \ll p_2.$$

Proof. We only show the first claim, the other proofs are similar. Let $p_1 = (k, \ell)$. Since $A_v^{T,i} \subseteq T^{\leq \lceil m/2 \rceil}$, p_1 has been added to $A_v^{T,k}$ while processing T^k in step k , and $k \leq i$. Clearly, from the way S is handled in lines 1–5 of Figure 3 (a), we have $p_1 \notin IP(T^k \cap M, A_j^{T,k-1})$, for $s^{TL,k-1} \leq j < v$. Since $s^{TL,k-1} \leq s^{TL,i-1} \leq u < v$, there is some $p_2 \in A_u^{T,k-1}$ such that $p_1 \gg p_2$ or $p_1 \ll p_2$. But the second case would imply $p_2 \in T^{k'}$ for some $k' > k$ which is impossible during the first k steps of our construction. Finally observe that the algorithm never removes matches while updating an antichain, thus p_2 is still present in $A_u^{T,i}$. \square

Lemma 3.2. The following holds:

$$a) \forall 1 \leq i \leq \lceil m/2 \rceil \forall v: v < s^{TL,i} \iff A_v^{T,i} \text{ or } A_v^{L,i} \text{ is TL-complete.}$$

$$b) \forall 1 \leq i \leq \lceil m/2 \rceil \forall v: v < s^{BR,i} \iff A_v^{B,i} \text{ or } A_v^{R,i} \text{ is BR-complete.}$$

Proof. We only prove the first claim, the other one is similar.

If. By contradiction, let i be the first step such that $A_v^{T,i}$ or $A_v^{L,i}$ is TL-complete, but $v \geq s^{TL,i}$. Clearly $v \neq s^{TL,i-1}$, otherwise the TL-completeness would have been detected by the algorithm shown in Figure 4 (a), and thus, contradicting the property of v , we would have $v < s^{TL,i} = s^{TL,i-1} + 1$. Hence $v > s^{TL,i-1}$. By the TL-completeness, there is some match $(k, \ell) \in A_v^{T,i} \cup A_v^{L,i}$ such that $1 \leq k, \ell \leq i$. Furthermore, by Lemma 3.1, there exists some match $(k', \ell') \in A_{v-1}^{T,i} \cup A_{v-1}^{L,i}$ such

that $(k', \ell') \ll (k, l)$. But then $1 \leq k', \ell' < i$, and therefore either $A_{v-1}^{T,i}$ or $A_{v-1}^{L,i}$ would be TL-complete after step $i-1$, a contradiction to the choice of i .

Only if. Obvious from the management of the start indices. \square

Lemma 3.3. For all i, u define $A_u^{TL,i} := A_u^{T,i} \cup A_u^{L,i}$ and $A_u^{BR,i} := A_u^{B,i} \cup A_u^{R,i}$. Then

- a) $\forall 0 \leq i \leq \lceil m/2 \rceil \forall 1 \leq u \leq \min\{e^{T,i}, e^{L,i}\}$: $A_u^{TL,i}$ is an antichain.
- b) $\forall 0 \leq i \leq \lceil m/2 \rceil \forall 1 \leq u \leq \min\{e^{B,i}, e^{R,i}\}$: $A_u^{BR,i}$ is an antichain.

Proof. We prove the first claim by induction on i . The base $i = 0$ is trivial because $A^{T,0} = A^{L,0} = \emptyset$. For the induction step $i-1 \rightarrow i$, we consider three different cases.

Case a: $1 \leq u < s^{TL,i-1}$. Then $A_u^{T,i} = A_u^{T,i-1}$ and $A_u^{L,i} = A_u^{L,i-1}$ (see lines 15 and 30 in Figure 3(a), respectively). Thus, by the induction hypothesis, $A_u^{TL,i}$ is an antichain.

Case b: $s^{TL,i-1} \leq u \leq \min\{e^{T,i-1}, e^{L,i-1}\}$. By definition the set $T := IP(S, A_u^{T,i-1})$ added to $A_u^{T,i}$ in line 3 (Figure 3(a)) is incomparable to $A_u^{T,i-1}$, but it is also incomparable to $A_u^{L,i}$ as we now demonstrate. Let $(k, \ell) \in IP(S, A_u^{T,i-1})$ and $(k', \ell') \in A_u^{L,i}$. Observe $k = i$ and $\ell \geq i$. Also note that $k' > \ell'$ and $\ell' \leq i$ because $A_u^{L,i} \subseteq L^{\leq i}$. Thus $(k, \ell) \ll (k', \ell')$ would contradict $\ell \geq i \geq \ell'$. Furthermore, $(k', \ell') \ll (k, \ell)$ would imply $\ell' < k' < k = i$, i.e., $A_u^{L,i-1}$ would be TL-complete, a contradiction to Lemma 3.2 and the choice of u . Similar arguments can be used for the set $L := IP(S, A_u^{L,i-1})$ added to $A_u^{L,i}$ in line 19. Finally note that $T \subseteq T^i$ and $L \subseteq L^i$ are also incomparable.

Case c: $\min\{e^{T,i-1}, e^{L,i-1}\} < u \leq \min\{e^{T,i}, e^{L,i}\}$. Clearly, this case is only possible if $u = e^{T,i} = e^{T,i-1} + 1$ or $u = e^{L,i} = e^{L,i-1} + 1$. If both conditions hold, then $A_u^{T,i} \subseteq T^i \cap M$ (lines 1 and 7) and $A_u^{L,i} \subseteq L^i \cap M$ (lines 17 and 23), thus their union obviously makes up an antichain. Otherwise, only one new antichain is generated whereas the other one is updated, and we can argue as in the second case to show that both antichains are incomparable.

The proof of the second claim is similar. \square

Lemma 3.4. Let $1 \leq i \leq \lceil m/2 \rceil$. Then the following holds:

- a) $\forall j \leq \max\{e^{T,i}, e^{L,i}\} \forall p_j \in A_j^{TL,i} \exists p_1 \in A_1^{TL,i}, \dots, p_{j-1} \in A_{j-1}^{TL,i}$:
 $p_1 \ll \dots \ll p_j$.
- b) $\forall j \leq \max\{e^{B,i}, e^{R,i}\} \forall p_j \in A_j^{BR,i} \exists p_1 \in A_1^{BR,i}, \dots, p_{j-1} \in A_{j-1}^{BR,i}$:
 $p_1 \gg \dots \gg p_j$.

Proof. We prove the first claim by choosing p_v for $v = j-1, \dots, 1$.

Consider step $j' \leq i$ when p_{v+1} was added to $A_{v+1}^{TL,j'} \subseteq A_{v+1}^{TL,i}$. Then Lemma 3.1 implies the existence of p_v if $v \geq s^{TL,j'-1}$. Otherwise, by Lemma 3.2, $A_v^{T,j'-1}$ or $A_v^{L,j'-1}$ has been detected to be TL-complete before step j' , i.e., $A_v^{TL,j'-1}$ contains a match (k', ℓ') such that $k', \ell' < j'$. But p_{v+1} is of the form (k, ℓ) with $k, \ell \geq j'$, thus we can choose $p_v := (k', \ell')$.

Similar arguments can be used for the second claim. \square

Lemma 3.5. For $0 \leq i \leq \lceil m/2 \rceil$, there are two chains

$$C^{TR,i}, C^{BL,i} \subseteq T^{\leq i} \cup L^{\leq i} \cup B^{\leq i} \cup R^{\leq i}$$

of length $e^{T,i} + e^{R,i}$ and $e^{B,i} + e^{L,i}$, respectively.

Proof. We prove the existence of the first chain $C^{TR,i}$ by induction on i . The base $i = 0$ is trivial. For the induction step $(i-1) \rightarrow i$, we have to analyse the situations which cause $e^{T,i} + e^{R,i}$ to be greater than $e^{T,i-1} + e^{R,i-1}$. One such situation is given in lines 7–14 of Figure 3(a) if the condition in line 10 is not satisfied because then $e := e^{T,i} = e^{T,i-1} + 1$ and $\tilde{e} := e^{R,i} = e^{R,i-1}$. But since $IP(A_e^{R,i-1}, A_e^{T,i}) \neq A_e^{R,i-1}$ there exist two comparable matches $c^T \in A_e^{T,i}$ and $c^R \in A_e^{R,i-1}$. More precisely, since $c^T \in T^i$ and $c^R \in R^{\leq i-1}$, we must have $(k, \ell) \ll (k', \ell')$. Thus, by Lemma 3.4, we can construct a chain

$$p_1 \ll \dots \ll p_{e-1} \ll c^T \ll c^R \ll p'_{\tilde{e}-1} \ll \dots \ll p'_1$$

of length $e + \tilde{e}$.

Similar arguments can be used for the remaining situations and for the other chain. \square

Our next task is to reveal the structure in D^{TR} and D^{BL} . We shall show that for each deleted match there always is some antichain which is incomparable to this match. In order to prove this property, we keep track of each deleted match by *assigning* it to some antichain during the construction process. More precisely, whenever an antichain A is removed due to the existence of some other antichain B which inactivates it, all matches in A are assigned to B , e.g., considering the situation in Figure 2(d), the match $(3, 8)$ is assigned to $A_3^{T,3}$. Furthermore, all previously deleted matches assigned to A now also belong to B . The assigned matches are inherited when an antichain is updated, e.g., in Figure 2(e), $(3, 8)$ also belongs to $A_3^{T,4}$. These rules guarantee that after step i , each deleted match is assigned to exactly one antichain in $A^{T,i} \cup A^{L,i} \cup A^{B,i} \cup A^{R,i}$. We write $D(A)$ to denote the set of matches assigned to an antichain A .

Lemma 3.6. Let $1 \leq i \leq \lceil m/2 \rceil$, and assume $(k, \ell) \in D(A)$ for some antichain A in $A^{T,i}$, $A^{L,i}$, $A^{B,i}$, or $A^{R,i}$. Then

$$\text{a) } (k, \ell) \in D^{TR} \implies \forall (k', \ell') \in A: k \leq k' \wedge \ell \geq \ell'.$$

$$\text{b) } (k, \ell) \in D^{BL} \implies \forall (k', \ell') \in A: k \geq k' \wedge \ell \leq \ell'.$$

Proof. For the first claim, let us assume (k, ℓ) was assigned to A while executing line 11 in Figure 3 (a) during step $j \leq i$ (the following arguments can analogously be applied to the other instructions which modify D^{TR}). Thus $A = A_e^{T,i}$, where $e = e^{T,j}$. Now we consider two cases concerning the status of (k, ℓ) before step j .

Case a: $(k, \ell) \in A_e^{R,j-1} \subseteq R^{\leq j-1}$, $\tilde{e} = e^{R,j-1}$. Then $\ell > n - j + 1$. From lines 1, 6, 7, and 10 we see that (k, ℓ) is incomparable to any match (k'', ℓ'') in $A_e^{T,j}$. But $A_e^{T,j} \subseteq T^j$, thus $k'' = j$ and $\ell'' \leq n - j + 1$. Hence, the incomparability implies $k \leq j$. Now observe that $A_e^{T,j}$ is the first constructed part of $A_e^{T,i}$, later extensions are taken from T^{j+1}, \dots, T^i . Thus every match $(k', \ell') \in A_e^{T,i}$ fulfills $k' \geq j$ and $\ell' \leq n - j + 1$, and the claim follows.

Case b: (k, ℓ) is assigned to $A_e^{R,j-1}$. We can inductively assume

$$\forall (k'', \ell'') \in A_e^{R,j-1}: k \leq k'' \wedge \ell \geq \ell''.$$

Deleted matches are never assigned to empty antichains. Thus there is at least one match $(k'', \ell'') \in A_e^{R,j-1}$, and we can prove as in the first case that $k'' \leq k'$ and $\ell'' \geq \ell'$. Hence we have $k \leq k'$ and $\ell \geq \ell'$.

The proof of the second claim is similar and therefore omitted. \square

Lemma 3.7. Let $1 \leq i \leq \lceil m/2 \rceil$. Then the following holds:

- a) $\forall 1 \leq u \leq e^{T,i}: D^{BL} \cap D(A_u^{T,i}) \neq \emptyset \implies A_u^{L,i} = \emptyset \wedge A_u^{T,i}$ is TL-complete.
- b) $\forall 1 \leq u \leq e^{L,i}: D^{TR} \cap D(A_u^{L,i}) \neq \emptyset \implies A_u^{T,i} = \emptyset \wedge A_u^{L,i}$ is TL-complete.
- c) $\forall 1 \leq u \leq e^{B,i}: D^{TR} \cap D(A_u^{B,i}) \neq \emptyset \implies A_u^{R,i} = \emptyset \wedge A_u^{B,i}$ is BR-complete.
- d) $\forall 1 \leq u \leq e^{R,i}: D^{BL} \cap D(A_u^{R,i}) \neq \emptyset \implies A_u^{B,i} = \emptyset \wedge A_u^{R,i}$ is BR-complete.

Proof. We again only show the first claim. From lines 10 and 11 in Figure 3 (a), we see that all matches assigned there to $A_u^{T,i}$ are either placed into D^{TR} , or they have been assigned before to some non-complete antichain in $A^{R,i-1}$. But concerning the latter case, we see from lines 26 and 27 in Figure 3 (b) that any such match has been put into D^{TR} as well, or again belongs to some non-complete antichain in $A^{T,j}$, $j < i$. Repeating this argument, we conclude that all matches assigned to $A^{T,i}$ are contained in D^{TR} . The only exception is given by lines 8 and 9 in Figure 4 (a), where deleted matches are assigned to $A_u^{T,i}$, but added to D^{BL} . But then, from lines 3, 4, and 13, the claim follows. \square

Lemma 3.8. All matches assigned to an antichain A are pairwise incomparable, thus by Lemma 3.6, they extend the antichain to a larger one.

Proof. Whenever a match is deleted, the algorithm always removes a complete antichain. By induction, this antichain B together with its assigned matches forms a larger antichain C . If there already is a set of matches D assigned to A (which is only possible when A is detected to be complete), then, following the arguments given in the proof of Lemma 3.7, $C \subseteq D^{BL}$ and $D \subseteq D^{TR}$ or vice versa, and Lemma 3.6 immediately implies that B and D are pairwise incomparable. \square

We are now prepared to construct a minimal decomposition of M . We start by decomposing $M \setminus (D^{TR} \cup D^{BL})$, the deleted matches are later considered in Theorem 3.9 below. The construction is as follows. Using Lemma 3.3, we combine the first $e^{TL} := \min\{e^{T, \lceil m/2 \rceil}, e^{L, \lceil m/2 \rceil}\}$ antichains in $A^{T, \lceil m/2 \rceil}$ and $A^{L, \lceil m/2 \rceil}$ to larger ones. We also connect the first $e^{BR} := \min\{e^{B, \lceil m/2 \rceil}, e^{R, \lceil m/2 \rceil}\}$ antichains in $A^{B, \lceil m/2 \rceil}$ to the corresponding ones in $A^{R, \lceil m/2 \rceil}$. For example, in Figure 2(e), we have $e^{T, \lceil m/2 \rceil} = e^{B, \lceil m/2 \rceil} = 3$ and $e^{L, \lceil m/2 \rceil} = e^{R, \lceil m/2 \rceil} = 2$, thus this generates four combined antichains. Concerning the remaining antichains we consider four different cases.

Case a: $e^{T, \lceil m/2 \rceil} \leq e^{L, \lceil m/2 \rceil}$ and $e^{B, \lceil m/2 \rceil} \geq e^{R, \lceil m/2 \rceil}$. Then we leave the remaining antichains as they are and have $p := e^{L, \lceil m/2 \rceil} + e^{B, \lceil m/2 \rceil}$ antichains in total. But by Lemma 3.5, there also exists a chain of this length. Thus, by Dilworth's theorem, the decomposition is minimal.

Case b: $e^{T, \lceil m/2 \rceil} > e^{L, \lceil m/2 \rceil}$ and $e^{B, \lceil m/2 \rceil} \leq e^{R, \lceil m/2 \rceil}$. Similar to the first case we have $p := e^{T, \lceil m/2 \rceil} + e^{R, \lceil m/2 \rceil}$ antichains, and also a chain of this length.

Case c: $e^{T, \lceil m/2 \rceil} \leq e^{L, \lceil m/2 \rceil}$ and $e^{B, \lceil m/2 \rceil} < e^{R, \lceil m/2 \rceil}$. From the management of the start and end indices, we have $e^{T, \lceil m/2 \rceil} \geq s^{TL, \lceil m/2 \rceil} - 1$. Thus, by Lemma 3.2, $A_u^{L, \lceil m/2 \rceil}$ is not TL-complete for $u > e^{T, \lceil m/2 \rceil}$. This implies $k > \lceil m/2 \rceil$ and $\ell \leq \lceil m/2 \rceil$ for any match $(k, \ell) \in A_u^{L, \lceil m/2 \rceil} \subseteq L^{\leq \lceil m/2 \rceil}$. For all $v > e^{B, \lceil m/2 \rceil}$ and $(k', \ell') \in A_v^{R, \lceil m/2 \rceil}$ we similarly have $k' \leq \lceil m/2 \rceil$ and $\ell' > n - \lfloor m/2 \rfloor \geq \lceil m/2 \rceil$. Thus $A_u^{L, \lceil m/2 \rceil}$ and $A_v^{R, \lceil m/2 \rceil}$ are incomparable. Now assume $e^{L, \lceil m/2 \rceil} \geq e^{R, \lceil m/2 \rceil}$. Then we can connect all remaining antichains in $A^{R, \lceil m/2 \rceil}$ to corresponding ones in $A^{L, \lceil m/2 \rceil}$ and obtain $p := e^{L, \lceil m/2 \rceil} + e^{B, \lceil m/2 \rceil}$ antichains in total, thus again a minimal decomposition. If $e^{L, \lceil m/2 \rceil} < e^{R, \lceil m/2 \rceil}$, then similarly $p := e^{T, \lceil m/2 \rceil} + e^{R, \lceil m/2 \rceil}$ is the optimal length of a chain in $M \setminus (D^{TR} \cup D^{BL})$.

Case d: $e^{T, \lceil m/2 \rceil} > e^{L, \lceil m/2 \rceil}$ and $e^{B, \lceil m/2 \rceil} > e^{R, \lceil m/2 \rceil}$. Finding a minimal decomposition is slightly more complicated in this case. Consider the following algorithm. Starting with $u := e^{T, \lceil m/2 \rceil}$ and $v := e^{R, \lceil m/2 \rceil} + 1$, we check whether $A_u^{T, \lceil m/2 \rceil}$ and $A_v^{B, \lceil m/2 \rceil}$ are incomparable. If they are not, then we backup u and v in \tilde{u} and \tilde{v} , respectively, and increase v by one. Otherwise the antichains are connected, u is set to $u - 1$, and v is set to $v + 1$. We repeat this until all remaining antichains in either $A^{T, \lceil m/2 \rceil}$ or $A^{B, \lceil m/2 \rceil}$ have been used, i.e., $u = e^{L, \lceil m/2 \rceil}$ or $v > e^{B, \lceil m/2 \rceil}$. Then the total number of antichains is $p := u + e^{B, \lceil m/2 \rceil}$. Thus, if $u = e^{L, \lceil m/2 \rceil}$, we have $p = e^{L, \lceil m/2 \rceil} + e^{B, \lceil m/2 \rceil}$, and the decomposition is optimal. Now assume $u > e^{L, \lceil m/2 \rceil}$.

If \tilde{u} and \tilde{v} are unused, then all remaining antichains in $A^{B, \lceil m/2 \rceil}$ have been connected to corresponding antichains in $A^{T, \lceil m/2 \rceil}$, and we have $p = e^{T, \lceil m/2 \rceil} + e^{R, \lceil m/2 \rceil}$. Hence, in this case the decomposition is also a minimal one. Finally assume that \tilde{u} and \tilde{v} have been used for saving u and v at least once. Then for $j = \tilde{v} + 1, \dots, e^{B, \lceil m/2 \rceil}$, $A_j^{B, \lceil m/2 \rceil}$ has been connected to $A_{\tilde{u} + \tilde{v} - j}^{T, \lceil m/2 \rceil}$, and we have $u = \tilde{u} - (e^{B, \lceil m/2 \rceil} - \tilde{v})$. Thus $p = \tilde{u} - (e^{B, \lceil m/2 \rceil} - \tilde{v}) + e^{B, \lceil m/2 \rceil} = \tilde{u} + \tilde{v}$. But from the properties of \tilde{u} and \tilde{v} , it can be shown (similar to the proof of Lemma 3.5) that there is a chain of length $\tilde{u} + \tilde{v}$ which contains two matches $p_1 \in A_{\tilde{u}}^{T, \lceil m/2 \rceil}$ and $p_2 \in A_{\tilde{v}}^{B, \lceil m/2 \rceil}$. Hence, the constructed decomposition is optimal.

Let us consider our example. Case d applies to the situation in Figure 2(e), and $A_3^{T,4}$ is compared with $A_3^{B,4}$. Since these antichains are incomparable, they are connected, and we obtain a decomposition consisting of 5 antichains in total.

Theorem 3.9. The length of an LCS in M equals p as defined in the four cases above.

Proof. Consider a combined antichain A of the decomposition. Assume an antichain $A_u^{T, \lceil m/2 \rceil} \in A^{T, \lceil m/2 \rceil}$ is one component of it (otherwise, we can handle the following construction in a similar way).

Case a: $A_u^{T, \lceil m/2 \rceil}$ is the only component of A . Then we extend A with the set B of deleted matches assigned to $A_u^{T, \lceil m/2 \rceil}$. Lemma 3.8 guarantees that the result is still an antichain.

Case b: $A_u^{T, \lceil m/2 \rceil}$ has been combined with $A_u^{L, \lceil m/2 \rceil}$. By Lemma 3.7, $B \subseteq D^{TR}$. Let $(k, \ell) \in A_u^{L, \lceil m/2 \rceil}$ and $(k', \ell') \in A_u^{T, \lceil m/2 \rceil}$. Now $(k, \ell) \in L^{\lceil m/2 \rceil}$, the incomparability of (k, ℓ) and (k', ℓ') , and $(k', \ell') \in T^{\lceil m/2 \rceil}$ imply that $k \geq k' \wedge \ell \leq \ell'$. Now consider a match $(k'', \ell'') \in B$. By Lemma 3.6, we have $k \geq k' \geq k''$ and $\ell \leq \ell' \leq \ell''$. Hence, $A_u^{L, \lceil m/2 \rceil}$ is incomparable to B . We can use a similar way to show that the set C of deleted matches assigned to $A_u^{L, \lceil m/2 \rceil}$ is a subset of D^{BL} and incomparable to $A_u^{T, \lceil m/2 \rceil}$. Finally, B and C are clearly incomparable as well. This implies that $A_u^{T, \lceil m/2 \rceil} \cup A_u^{L, \lceil m/2 \rceil} \cup B \cup C$ is still an antichain.

Case c: $A_u^{T, \lceil m/2 \rceil}$ has been combined with some other antichain $D \in A^{B, i}$. Then, similar to the proof of the second case, we can show that the union of A and the two corresponding sets of assigned matches still make up an antichain.

By handling each combined antichain in this way, we can construct a decomposition of M without generating any additional antichains. The proof is complete. \square

Figure 2(f) illustrates the corresponding decomposition for our example.

4. IMPLEMENTATION

We now describe an efficient implementation for the given algorithm and analyse its time and space complexity.

All new antichains created in step i are extensions from antichains generated during step $i-1$. Furthermore, the only antichains used for decomposing M are from

the last step. Thus for the implementation it is sufficient to update the antichains of interest. The same is true for the start and end indices, and we thus sometimes drop the index i from now on. The necessary information for each actual antichain can be kept in one single number as follows. Let $1 \leq i \leq \lceil m/2 \rceil$ and $1 \leq u \leq e^{T,i}$. We define $ThreshT[u]$ as the leftmost column used by some match in $A_u^{T,i}$, i.e.,

$$ThreshT[u] := \min\{\ell \mid \exists k: (k, \ell) \in A_u^{T,i}\}.$$

For example, in Figure 2 (b), $ThreshT[1] = 3$, and in Figure 2 (d), $TopThresh[1] = 2$, $ThreshT[2] = 3$, and $ThreshT[3] = 6$. To update this array in each step, we use an auxiliary array $LeftPos$ on $\Sigma \times [1 : n + 1]$ given by

$$LeftPos[c, \ell] := \min(\{n + 1\} \cup \{j \mid \ell \leq j \leq n \wedge y_j = c\}),$$

i.e., $LeftPos[a_i, \ell]$ equals the column number of the leftmost occurrence of a match in row i located right to column ℓ , and equals $n + 1$ if there is no such match. In our example ($y = cbabbacac$), we obtain the following values:

a	3	3	3	6	6	6	8	8	10	10
b	2	2	4	4	5	10	10	10	10	10
c	1	7	7	7	7	7	7	9	9	10

Now it is not difficult to see that the following routine correctly updates $ThreshT$ when processing T^i , representing lines 1–7 in Figure 3 (a). (Similar procedures are used in [4, 29, 30] to determine *contours* which correspond to the antichains used here.)

```

 $k := LeftPos[a_i, i];$ 
For  $u := s^{TL}$  To  $e^T$  Do {
   $j := ThreshT[u];$ 
  If  $k \leq j$  And  $k \leq n - i + 1$  Then {
     $ThreshT[u] := k; k := LeftPos[a_i, j + 1];$ 
  };
};
If  $k \leq n - i + 1$  Then {  $e^T := e^T + 1; ThreshT[e^T] := k$  };

```

For $A^{L,i}$, $A^{B,i}$, and $A^{R,i}$ we introduce additional arrays $ThreshL$, $ThreshB$, and $ThreshR$ which similarly store the topmost rows, rightmost columns, and bottommost rows used by the corresponding antichains. To handle them analogously to $ThreshT$, we also need three more auxiliary arrays given by

$$\begin{aligned}
TopPos[c, k] &:= \min(\{m + 1\} \cup \{j \mid k \leq j \leq m \wedge x_j = c\}), & (1 \leq k \leq m + 1), \\
RightPos[c, \ell] &:= \max(\{0\} \cup \{j \mid 1 \leq j \leq \ell \wedge y_j = c\}), & (0 \leq \ell \leq n), \\
BottomPos[c, k] &:= \max(\{0\} \cup \{j \mid 1 \leq j \leq k \wedge x_j = c\}), & (0 \leq k \leq m).
\end{aligned}$$

Note that in Figure 3 and Figure 4, each test for the incomparability of two antichains can be replaced by a rather simple conditional statement. For example, considering line 10 in Figure 3 (a), we know that all matches in T^i are located to the left of any match in $R^{\leq i-1}$. Thus, with $e := e^{T,i}$ and $\tilde{e} := e^{R,i}$, A_e^T and $A_{\tilde{e}}^R$ are incomparable if

and only if A_e^R is also completely contained in the first i rows, i.e., $\text{ThreshR}[\tilde{e}] \leq i$. The algorithm presented in Figure 6 shows how the other situations are handled. It also makes use of some special implementation details which cannot be discussed here, e.g., the construction starts with the bottommost row instead of the topmost one when m is even. In Figure 6 some lines are marked with a dot (\bullet) on their left sides. These lines are used for the construction of an LCS and should be ignored for the moment.

The complexity of the algorithm may be deduced as follows. The four auxiliary arrays can be easily preprocessed in $O(ns)$ time and space, where $s = |\Sigma|$. Clearly, during one of the $\lceil m/2 \rceil$ iterations of the main loop, none of the four inner **While**-loops takes more than $O(p)$ time, and when determining p , at most $\lceil m/2 \rceil$ pairs of antichains have to be compared. Thus the algorithm takes at most $O(ns + mp)$ time. Furthermore, observe that the j -th antichain in A^T (which is added to A^T during some step $i \geq j$) must contain a match (k, ℓ) with $\ell \leq n - (p - j)$, otherwise it would be impossible to construct a chain of length p . But then this antichain is detected to be TL-complete after step $n - (p - j)$, therefore it is only considered for at most $n - (p - j) - i \leq n - p$ times in the corresponding **While**-loop (lines 59–65). Similar arguments can be given for antichains in A^L , A^B , and A^R . Hence, we have shown the following theorem.

Theorem 4.1. Let $x, y \in \Sigma^+$, $m = |x|$, $n = |y|$, $m \leq n$, and $s = |\Sigma|$. Then the length p of an LCS of x and y can be computed in $O(ns + \min\{mp, p(n - p)\})$ time and $O(ns)$ space.

This result has been achieved before by Rick [29, 30], and in fact, the algorithm presented here is some kind of dualization of Rick's method, but our algorithm is significantly faster as we shall show in Section 6.

5. CONSTRUCTION OF AN LCS IN LINEAR SPACE

This section deals with the generation of an LCS. The idea is to apply the divide-and-conquer scheme [5, 16, 20] which first identifies at least one point of an LCS such that this LCS is splitted into two parts of roughly the same size. Then the remainder is computed by recursive calls. The method presented here usually determines two LCS-neighbouring matches c^{TL} and c^{BR} which are located in $T^{\leq \lceil m/2 \rceil} \cup L^{\leq \lceil m/2 \rceil}$ and $B^{\leq \lceil m/2 \rceil} \cup R^{\leq \lceil m/2 \rceil}$, respectively. This is accomplished as follows.

In each step i of the construction described in Section 2, we subsequently update the following variables:

- p^{TL} is the match which caused $A_s^{T,i}$ or $A_s^{L,i}$ to become TL-complete, where $s = s^{TL,i} - 1$. For example, in Figure 2 (c), $p^{TL} = (2, 2)$, and in Figure 2 (d) and (e), $p^{TL} = (3, 3)$.
- p^{BR} has a corresponding meaning for the last BR-complete antichain in $A^{B,i}$ and $A^{R,i}$, e.g., in Figure 2 (d), $p^{BR} = (6, 7)$.


```

Determine TopPos and LeftPos;
Determine BottomPos and RightPos;
For  $u := 0$  To  $\lceil m/2 \rceil$  Do {
  ThreshT[u] := 0; ThreshL[u] := 0;
5 };
For  $u := 0$  To  $\lfloor m/2 \rfloor$  Do {
  ThreshB[u] :=  $n + 1$ ; ThreshR[u] :=  $m + 1$ ;
};
 $t := 1$ ;  $\ell := 1$ ;  $b := m$ ;  $r := n$ ;
10  $s^{TL} := 1$ ;  $e^T := 0$ ;  $e^L := 0$ ;
 $s^{BR} := 1$ ;  $e^B := 0$ ;  $e^R := 0$ ;

If Odd( $m$ ) Then Goto Line 57;

While  $t \leq b$  Do {      (* Main loop *)
   $k := \text{RightPos}[x_b, r]$ ;      (* Update  $A^B$  *)
   $u := s^{BR}$ ;
15 While  $u \leq e^B$  Do {
     $j := \text{ThreshB}[u]$ ;
    If  $k \geq j$  Then {
      ThreshB[u] :=  $k$ ;  $k := \text{RightPos}[x_b, j - 1]$ ;
20 };
     $u := u + 1$ ;
  };
  If  $k \geq \ell$  Then {
     $e^B := u$ ; ThreshB[eB] :=  $k$ ;
25 If ThreshL[eL]  $\geq b$  Then  $e^L := e^L - 1$ 
    • Else Update  $c^B, c^L, \ell^{BL}$ ;
  };
   $k := \text{BottomPos}[y_r, b - 1]$ ;      (* Update  $A^R$  *)
   $u := s^{BR}$ ;
30 While  $u \leq e^R$  Do {
     $j := \text{ThreshR}[u]$ ;
    If  $k \geq j$  Then {
      ThreshR[u] :=  $k$ ;  $k := \text{BottomPos}[y_r, j - 1]$ ;
35 };
     $u := u + 1$ ;
  };
  If  $k \geq t$  Then {
     $e^R := u$ ; ThreshR[eR] :=  $k$ ;
    If ThreshT[eT]  $\geq r$  Then  $e^T := e^T - 1$ 
    • Else Update  $c^T, c^R, \ell^{TR}$ ;
  };
  (* Check for BR-complete antichains *)
  If ThreshB[sBR] =  $r$  Then {
    If  $s^{BR} > e^R$  Then {
45 If ThreshT[eT]  $\geq r$  Then  $e^T := e^T - 1$ 
    • Else Update  $c^T, c^R, \ell^{TR}$ ;
    };
     $s^{BR} := s^{BR} + 1$ ;
  } Else If ThreshR[sBR] =  $b$  Then {
    If  $s^{BR} > e^B$  Then {
50 If ThreshL[eL]  $\geq b$  Then  $e^L := e^L - 1$ 
    • Else Update  $c^B, c^L, \ell^{BL}$ ;
    };
     $s^{BR} := s^{BR} + 1$ ;
55 };
   $t := t + 1$ ;  $\ell := \ell + 1$ ;

   $k := \text{LeftPos}[x_t, \ell]$ ;      (* Update  $A^T$  *)
   $u := s^{TL}$ ;
  While  $u \leq e^T$  Do {
     $j := \text{ThreshT}[u]$ ;
    If  $k \leq j$  Then {
      ThreshT[u] :=  $k$ ;  $k := \text{LeftPos}[x_t, j + 1]$ ;
60 };
     $u := u + 1$ ;
  };
  If  $k \leq r$  Then {
     $e^T := u$ ; ThreshT[eT] :=  $k$ ;
    If ThreshR[eR]  $\leq t$  Then  $e^R := e^R - 1$ 
    • Else Update  $c^T, c^R, \ell^{TR}$ ;
  };
70 };
   $k := \text{TopPos}[y_t, t]$ ;      (* Update  $A^L$  *)
   $u := s^{TL}$ ;
  While  $u \leq e^L$  Do {
     $j := \text{ThreshL}[u]$ ;
    If  $k \leq j$  Then {
75 ThreshL[u] :=  $k$ ;  $k := \text{TopPos}[y_t, j + 1]$ ;
    };
     $u := u + 1$ ;
  };
  If  $k \leq b$  Then {
     $e^L := u$ ; ThreshL[eL] :=  $k$ ;
    If ThreshB[eB]  $\leq \ell$  Then  $e^B := e^B - 1$ 
    • Else Update  $c^B, c^L, \ell^{BL}$ ;
  };
80 };
  (* Check for TL-complete antichains *)
  If ThreshT[sTL] =  $\ell$  Then {
    If  $s^{TL} > e^L$  Then {
      If ThreshB[eB]  $\leq \ell$  Then  $e^B := e^B - 1$ 
      • Else Update  $c^B, c^L, \ell^{BL}$ ;
    };
     $s^{TL} := s^{TL} + 1$ ;
  } Else If ThreshL[sTL] =  $t$  Then {
    If  $s^{TL} > e^T$  Then {
      If ThreshR[eR]  $\leq t$  Then  $e^R := e^R - 1$ 
      • Else Update  $c^T, c^R, \ell^{TR}$ ;
    };
     $s^{TL} := s^{TL} + 1$ ;
  };
   $b := b - 1$ ;  $r := r - 1$ ;
100 };

(* Determine length  $p$  of an LCS *)
If  $e^T > e^L$  And  $e^B > e^R$  Then {
  If  $s^{TL} \leq e^L$  Then  $s^{TL} := e^L + 1$ ;
  If  $s^{BR} \leq e^R$  Then  $s^{BR} := e^R + 1$ ;
105  $u := e^T$ ;  $v := s^{BR}$ ;
  While  $u \geq s^{TL}$  And  $v \leq e^B$  Do {
    If ThreshT[u]  $\geq \text{ThreshB}[v]$ 
    Then  $u := u - 1$ 
    • Else {  $\tilde{u} := u$ ;  $\tilde{v} := v$  };
     $v := v + 1$ ;
  };
   $p := u + e^B$ ;
113 } Else  $p := \max\{e^L + e^B, e^T + e^R\}$ ;

```

Fig. 6. The $O(ns + \min\{mp, p(n-p)\})$ algorithm for determining the length p of an LCS.

- c^T and c^R are the two matches introduced in the proof of Lemma 3.5. They both lie in $C^{TR,i}$ and are neighbours in this chain. Furthermore, c^T and c^R are always located in the first i topmost rows and i rightmost columns of M , respectively.
- c^B and c^L have analogous properties for $C^{BL,i}$.
- ℓ^{TR} and ℓ^{BL} is the position of c^T in $C^{TR,i}$ and of c^L in $C^{BL,i}$, respectively. Also, $\ell^{TR} + 1$ and $\ell^{BL} + 1$ is the position of c^R in $C^{TR,i}$ and of c^B in $C^{BL,i}$, respectively.

Variables p^{TL} and p^{BR} can be easily updated. For example, consider lines 85–98 in Figure 6 where new TL-complete antichains are handled. Let $p^{TL} = (u, v)$. If the condition in line 86 is satisfied, then we know p^{TL} has to be set to the bottommost match located in the first t rows and column ℓ . Therefore two additional statements can be inserted between lines 86 and 87 such that u is set to $BottomPos[y_\ell, t]$ and v is set to ℓ . Similar statements apply for the situation in lines 92–98, and this completes the description of the management for p^{TL} . p^{BR} can be handled in a similar way.

c^T , c^R , and ℓ^{TR} must be updated whenever the length of $C^{TR,i}$ increases. These situations are indicated in lines 40, 46, 69, and 95 in Figure 6, and here we only sketch how to manage them. By arguments analogous to the ones given in the proof of Lemma 3.4, we have to distinguish two cases when updating c^T . If $s^{TL,i} > e^{T,i}$, then c^T is set to p^{TL} , otherwise c^T can be determined by some additional statements which are similar to the ones used for updating p^{TL} . In either case, we set ℓ^{TR} to $e^{T,i}$ because $e^{T,i}$ is the position of c^T in $C^{TR,i}$, as seen in the proof of Lemma 3.5. The management of c^B , c^L , and ℓ^{BL} is similar.

Now let us review the construction of the final decomposition given in the end of Section 3. If p is set to $e^{T,\lceil m/2 \rceil} + e^{R,\lceil m/2 \rceil}$, then we can use c^T and c^R as the appropriate matches for c^{TL} and c^{BR} . Similarly, if $p = e^{B,\lceil m/2 \rceil} + e^{L,\lceil m/2 \rceil}$, we establish $c^{TL} = c^L$ and $c^{BR} = c^B$. Finally, if a longest chain is determined by the algorithm described in case d of the construction (corresponding to lines 103–112 in Figure 6), and p is not set to one of the above values, then we can use the backup values \tilde{u} and \tilde{v} to determine $c^{TL} := (BottomPos[y_{\tilde{u}}, b], y_{\tilde{u}})$ and $c^{BR} := (TopPos[y_{\tilde{v}}, t], y_{\tilde{v}})$, where $\tilde{u} := ThreshT[\tilde{u}]$ and $\tilde{v} := ThreshB[\tilde{v}]$.

Before recursively calling the algorithm for the remaining parts of the LCS, we see it is necessary for our routine to not only work on the complete matrix of size $[1 : m] \times [1 : n]$, but also on any subarea $[k_1 : k_2] \times [\ell_1 : \ell_2]$. The necessary changes are quite straightforward, and we do not provide any details here. Moreover, it might be impossible to locate both c^{TL} and c^{BR} (e.g., when $|M| = 1$), but then one recursive call can simply be skipped.

Theorem 5.1. An LCS can be constructed in $O(ns + \min\{mp, p(n - p)\})$ time and $O(ns)$ space.

Proof. Clearly, for the top-level call, the additional overhead needed to keep track of the new variables is bounded by $O(m)$. Thus, not taking into account the time consumed by preprocessing or any recursive calls, we can assume the number

of elementary operations to be bounded by $d(m + \min\{mp, p(n - p)\})$, for some appropriate constant d . We first examine the bound $d(m + mp)$. Let $c^{TL} = (k, \ell)$ and $c^{BR} = (k', \ell')$ (if only one match has been determined, the analysis is similar). Consider the two first-level recursive calls concerning the areas M_1 and M_2 , where $M_1 := [1 : k - 1] \times [1 : \ell - 1]$ and $M_2 := [k' + 1 : m] \times [\ell' + 1 : n]$. Let p_1 and p_2 denote the length of an LCS in M_1 and M_2 , respectively, i. e., $p_1 + p_2 = p - 2$. Recall that c^{TL} is located in the first $\lceil m/2 \rceil$ rows and columns, i. e., the length of one side of M_1 is bounded by $\lceil m/2 \rceil - 1$. The same is true for M_2 , and thus the number of operations taken for both first-level calls is bounded by

$$d(\lceil m/2 \rceil - 1)(p_1 + 1) + d(\lceil m/2 \rceil - 1)(p_2 + 1) \leq dp \frac{m}{2}.$$

Repeating this argument, we obtain a $dmp/2^i$ bound for the at most 2^i i th-level recursive calls. Since recursion ends at level $\lceil \log(m/2) \rceil$, this sums up to at most $2 \cdot dmp$ for the complete algorithm.

Similar (but somewhat more complicated) arguments can be used to show the other bound $d(m + p(n - p))$. We refer to [15] for details.

Now finally observe that when comparing the divide-and-conquer routine with the algorithm which determines the length p of an LCS, we only need $O(\log m)$ additional stack space, and thus the $O(ns)$ space bound is still valid. \square

6. EXPERIMENTAL RESULTS

We compared our routine with the algorithm proposed by Rick [29, 30] which clearly outperforms any other method when constructing longest common subsequences of intermediate lengths. Rick's algorithm is also a flexible one, being very efficient for short and long LCS as well. It uses a strategy similar to the one presented here, but only constructs antichains (or *contours*) from the top and left side of M . While this substantially simplifies the implementation and also the preprocessing phase (i. e., we only have to compute *LeftPos* and *TopPos*), there are two severe drawbacks. First, in order to recover an LCS after determining its length, the so-called *dominant matches* must be saved during the construction of the contours, and this might take $\Omega(mn)$ space. Second, the number of checks of *Thresh*-values is significantly increased when decomposing M from only two sides. For an alphabet of size 8, Table 1 shows some sample results when determining p for different settings of m , n , and p .

The corresponding running times are presented in Table 2. Both algorithms were programmed in a straightforward way, using no special optimizations, and were tested on an Intel Pentium II at 300 MHz. It can be seen that our algorithm only takes about 70 % of the time needed by Rick's method when computing the length of an LCS which is of intermediate length. For very short or very long LCS our method slightly suffers from the additional overhead during the preprocessing phase, but is still very efficient.

Finally, we checked the running times and the consumed space when generating an LCS. Table 3 shows that in spite of the linear space restriction, our algorithm

Table 1. Frequency of checks of *Thresh*-values.

m	n	p	Rick [30]	New method
500	500	100	16864	14983
500	500	200	28962	23078
500	500	300	33276	23394
500	500	400	20384	13276
1500	1500	300	145129	126796
1500	1500	600	265107	216845
1500	1500	900	280026	207000
1500	1500	1200	172846	121516

Table 2. Running times in microseconds for determining the length p of an LCS.

m	n	p	Rick [30]	New method
500	500	100	3352	3626
500	500	200	5659	4725
500	500	300	6978	4890
500	500	400	5000	3516
1500	1500	300	24451	21868
1500	1500	600	46099	34835
1500	1500	900	54176	33791
1500	1500	1200	38791	22308

sometimes runs more than twice as fast as Rick's method. This is due to the significant overhead in Rick's routine which is caused by the additional statements responsible for saving the contours in memory.

Table 3. Running times in microseconds for constructing an LCS of length p .

m	n	p	Rick [30]	New method
500	500	100	6319	6044
500	500	200	14341	9066
500	500	300	19505	9890
500	500	400	15769	7802
750	750	250	23132	16374
750	750	400	39835	20495
750	750	550	38516	16758
750	750	700	16319	9945

Table 4. Allocated space in bytes
for constructing an LCS of length p .

m	n	p	Rick [30]	New method
500	500	100	64284	34072
500	500	200	143820	34072
500	500	300	199464	34072
500	500	400	176328	34072
750	750	250	219244	51072
750	750	400	390172	51072
750	750	550	396136	51072
750	750	700	193780	51072

7. CONCLUSION

We have investigated a new algorithm for the Longest Common Subsequence Problem. In spite of the quite complicated technical details necessary for the construction and analysis, the final routines proved to be very practical. More precisely, we have shown three results. First, we have presented a new fast method for determining the length of an LCS. Second, we have developed a linear space algorithm for constructing an LCS in $O(ns + \min\{mp, p(n-p)\})$ time, thus solving a previously open problem. And third, we have shown by some experimental results that this algorithm seems to be well-suited for many usual applications.

The presented method can be extended to find *all* LCS of two given strings while preserving the time complexity $O(ns + \min\{mp, p(n-p)\})$, which is the same time complexity as for Rick's algorithm. Details can be found in [15].

ACKNOWLEDGEMENT

We would like to thank Dr. F. Kurth and the anonymous referees for helpful comments.

(Received May 12, 2000.)

REFERENCES

-
- [1] A. V. Aho, D. S. Hirschberg, and J. D. Ullman: Bounds on the complexity of the longest common subsequence problem. *J. Assoc. Comput. Mach.* *23* (1976), 1, 1–12.
 - [2] A. Apostolico: Improving the worst-case performance of the Hunt–Szymanski strategy for the longest common subsequence of two strings. *Inform. Process. Lett.* *23* (1986), 63–69.
 - [3] A. Apostolico: Remarks on the Hsu–Du new algorithm for the longest common subsequence problem. *Inform. Process. Lett.* *25* (1987), 235–236.
 - [4] A. Apostolico and G. Guerra: The longest common subsequence problem revisited. *Algorithmica* *2* (1987), 315–336.
 - [5] A. Apostolico, S. Browne, and C. Guerra: Fast linear-space computations of longest common subsequences. *Theoret. Comput. Sci.* *92* (1992), 3–17.

- [6] F. Y. L. Chin and C. K. Poon: A fast algorithm for computing longest common subsequences of small alphabet size. *J. Inform. Process.* *13* (1990), 4, 463–469.
- [7] V. Chvátal and D. Sankoff: Longest common subsequences of two random strings. *J. Appl. Probab.* *12* (1975), 306–315.
- [8] V. Dančák and M. Paterson: Upper bounds for the expected length of a longest common subsequence of two binary sequences. In: *Proceedings 11th Annual Symp. on Theoretical Aspects of Computer Science (Lecture Notes in Computer Science 775)*, Springer-Verlag, Berlin 1994, pp. 669–678.
- [9] M. O. Dayhoff: Computer aids to protein sequence determination. *J. Theoret. Biol.* *8* (1965), 97–112.
- [10] M. O. Dayhoff: Computer analysis of protein evolution. *Sci. Amer.* *221* (1969), 1, 86–95.
- [11] J. G. Deken: Some limit results for longest common subsequences. *Discrete Math.* *26* (1979), 17–31.
- [12] R. P. Dilworth: A decomposition theorem for partially ordered sets. *Ann. of Math.* *51* (1950), 161–166.
- [13] K. S. Fu and B. K. Bhargava: Tree systems for syntactic pattern recognition. *IEEE Trans. Comput. C-22* (1973), 12, 1087–1099.
- [14] J. Gallant, D. Maier, and J. A. Storer: On finding minimal length superstrings. *J. Comput. System Sci.* *20* (1980), 50–58.
- [15] H. Goeman: Time and Space Efficient Algorithms for Decomposing Certain Partially Ordered Sets. PhD Thesis, Department of Computer Science, University of Bonn 1999. To appear in *Bayreuther Mathematische Schriften*.
- [16] D. S. Hirschberg: A linear space algorithm for computing maximal common subsequences. *Comm. ACM* *18* (1975), 6, 341–343.
- [17] D. S. Hirschberg: Algorithms for the longest common subsequence problem. *J. Assoc. Comput. Mach.* *24* (1977), 4, 664–675.
- [18] W. J. Hsu and M. W. Du: New algorithms for the LCS problem. *J. Comput. System Sci.* *29* (1984), 133–152.
- [19] J. W. Hunt and T. G. Szymanski: A fast algorithm for computing longest common subsequences. *Comm. ACM* *20* (1977), 5, 350–353.
- [20] S. K. Kumar and C. P. Rangan: A linear space algorithm for the LCS problem. *Acta Inform.* *24* (1987), 353–363.
- [21] R. Lowrance and R. A. Wagner: An extension of the string-to-string correction problem. *J. Assoc. Comput. Mach.* *22*, (1975), 2, 177–183.
- [22] S. Y. Lu and K. S. Fu: A sentence-to-sentence clustering procedure for pattern analysis. *IEEE Trans. Systems Man Cybernet. SMC-8*, (1978), 5, 381–389.
- [23] D. Maier: The complexity of some problem on subsequences and supersequences. *J. Assoc. Comput. Mach.* *25* (1978), 2, 322–336.
- [24] W. J. Masek and M. S. Paterson: A faster algorithm for computing string edit distances. *J. Comput. System Sci.* *20* (1980), 1, 18–31.
- [25] E. W. Myers: An $O(ND)$ difference algorithm and its variations. *Algorithmica* *1* (1986), 251–266.
- [26] N. Nakatsu, Y. Kambayashi, and S. Yajima: A longest common subsequence algorithm suitable for similar text strings. *Acta Inform.* *18* (1982), 171–179.
- [27] S. B. Needleman and C. S. Wunsch: A general method applicable to the search for similarities in the amino acid sequence of two proteins. *J. Molecular Biol.* *48* (1970), 443–453.
- [28] M. Paterson and V. Dančák: Longest common subsequences. In: *Proceedings, 19th Intern. Symp. on Mathematical Foundations of Computer Science (Lecture Notes in Computer Science 841)*, Springer Verlag, Berlin 1994, pp. 127–142.

- [29] C. Rick: New Algorithms for the Longest Common Subsequence Problem. Research Report No. 85123-CS, Department of Computer Science, University of Bonn 1994.
- [30] C. Rick: A new flexible algorithm for the longest common subsequence problem. In: Proceedings, 6th Annual Symp. on Combinatorial Pattern Matching (Lecture Notes in Computer Science 937), Springer Verlag, Berlin 1995, pp. 340–351.
- [31] D. Sankoff and R. J. Cedergren: A test for nucleotide sequence homology. *J. Molecular Biol.* *77* (1973), 159–164.
- [32] D. Sankoff and J. B. Kruskal: Time Warps, String Edits, and Macromolecules: The Theory And Practice of Sequence Comparison. Addison–Wesley, Reading, MA 1983.
- [33] E. Ukkonen: Algorithms for approximate string matching. *Inform. and Control* *64* (1985), 100–118.
- [34] R. A. Wagner: On the complexity of the extended string-to-string correction problem. In: Proceedings, 7th Ann. ACM Sympos. on Theory of Comput. 1975, pp. 218–223.
- [35] R. A. Wagner and M. J. Fischer: The string-to-string correction problem. *J. Assoc. Comput. Mach.* *21* (1974), 1, 168–173.
- [36] C. K. Wong and A. K. Chandra: Bounds for the string editing problem. *J. Assoc. Comput. Mach.* *28* (1976), 1, 13–18.
- [37] S. Wu, U. Manber, G. Myers, and W. Miller: An $O(NP)$ sequence comparison algorithm. *Inform. Process. Lett.* *35* (1990), 317–323.

Dr. Heiko Goeman and Prof. Dr. Michael Clausen, University of Bonn, Computer Science Department III, D-53117 Bonn. Germany.
e-mails: Heiko.Goeman@tlc.de, clausen@cs.uni-bonn.de