

Directed Acyclic Subsequence Graph¹

Zdeněk Troníček and Bořivoj Melichar

Department of Computer Science and Engineering
Czech Technical University
Karlovo náměstí 13, 121 35 Prague 2, Czech Republic
phone: ++420 2 2435 7287, fax: ++420 2 298098

e-mail: {tronicek,melichar}@fel.cvut.cz

Abstract. Directed Acyclic Subsequence Graph is an automaton, which accepts all subsequences of the given string. We introduce a left-to-right algorithm for incremental construction of DASG. The algorithm requires $\mathcal{O}(z)$ extra space and $\mathcal{O}(nz \log z)$ time for arbitrary alphabet ($\mathcal{O}(nz)$ for fixed alphabet), where $z = \min(|\Sigma|, n)$. The number of transitions can be reduced by encoding input symbols using k digits, where $k < \min(|\Sigma|, n)$. We introduce a left-to-right algorithm for incremental construction of DASG for $k = 2$. We show the extension of the algorithm for the set of strings and its application for the longest common subsequence problem.

Key words: Directed Acyclic Subsequence Graph, finite automaton, searching subsequences

1 Introduction

A subsequence of a string is any string obtained by deleting zero or more symbols from the given string. Directed Acyclic Subsequence Graph (DASG) is an automaton, which accepts all subsequences of the given text. It was introduced in [2] (preliminary version was published in [1]). DASG is analogous to Directed Acyclic Word Graph (DAWG) [3] using subsequences instead of substrings.

Let us suppose an alphabet Σ and a text $T = t_1 t_2 \dots t_n$ over this alphabet. DASG for the text T is an automaton $\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$, where Q is a finite set of states, Σ is an input alphabet, δ is a transition function, q_0 is the initial state and F is a set of final states. States are denoted by numbers in this article.

In [2], there is described a right-to-left algorithm for construction of DASG and encoding for reducing the number of transitions. In section 3 we introduce an incremental left-to-right algorithm for construction of DASG, and in section 4 its modification for encoded DASG. In section 5 we show the extension of the algorithm for a set of strings and its application for the longest common subsequence problem.

¹This research has been supported by GAČR grant No. 201/98/1155

2 Motivation

Let $Sub(T)$ denotes the set of all subsequences of the text $T = t_1t_2 \dots t_n$. The set $Sub(T)$ can be described recursively by the regular expression (ε is empty subsequence):

$$Sub_0 = \varepsilon$$

$$Sub_i = Sub_{i-1}(\varepsilon + t_i)$$

$$Sub(T) = Sub_n$$

For the set Sub_n then holds:

$$Sub_n = Sub_{n-1}(\varepsilon + t_n) = Sub_{n-2}(\varepsilon + t_{n-1})(\varepsilon + t_n) = (\varepsilon + t_1)(\varepsilon + t_2) \dots (\varepsilon + t_n)$$

This expression gives us the direction for construction of the nondeterministic version of DASG. The example of such nondeterministic finite automaton (NFA) is in Fig. 1. It also holds:

$$Sub_n = Sub_{n-1} + Sub_{n-1}t_n = \varepsilon + Sub_0t_1 + Sub_1t_2 + \dots + Sub_{n-2}t_{n-1} + Sub_{n-1}t_n$$

The last expression can be used for construction of the nondeterministic DASG without ε -transitions (the example is in Fig. 2). If all the symbols of T are different, we obtain directly the deterministic finite automaton (DFA). The example of deterministic DASG is in Fig. 3.

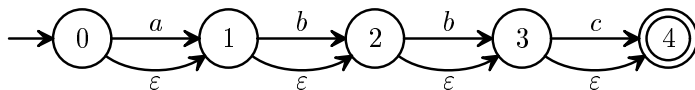


Figure 1: NFA accepting all subsequences of $T = abc$ (with ε -transitions).

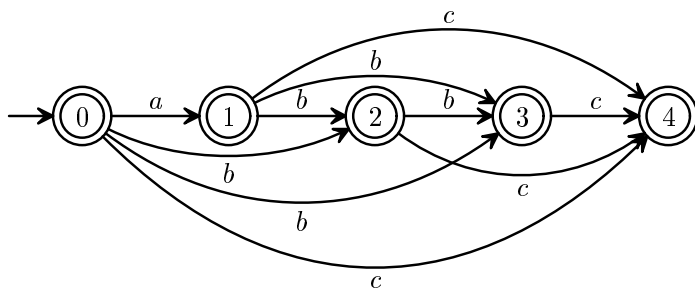
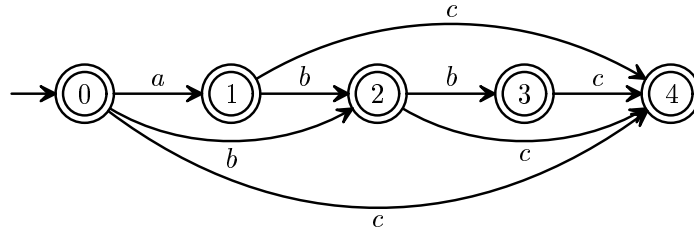


Figure 2: NFA accepting all subsequences of $T = abc$ (without ε -transitions).

3 Construction of DASG

Let us suppose an alphabet Σ and a text $T = t_1t_2 \dots t_n$ over this alphabet. For each symbol a of the alphabet Σ we will maintain the value f_a , which is the smallest number of the state not having an output transition labeled with a . We start with an automaton with the only state 0. Each state of the automaton is final.

Lemma 1: The automaton constructed by Algorithm 1 has $(n + 1)$ states.


 Figure 3: DASG for the string $T = abc$.

```

1: for each  $a \in \Sigma$  do
2:    $f_a \leftarrow 0$ 
3: end for
4: for  $k \leftarrow 1$  to  $n$  do
5:   add state  $k$ 
6:   for  $s \leftarrow f_{t_k}$  to  $(k-1)$  do
7:     add a transition labeled  $t_k$  from the state  $s$  to the state  $k$ 
8:   end for
9:    $f_{t_k} \leftarrow k$ 
10: end for
    
```

Figure 4: Algorithm 1 (incremental construction of DASG)

Proof: We start with the automaton with one state. The main cycle of the algorithm is performed n times. During each step of the cycle we add just one new state.

Lemma 2: The automaton constructed by Algorithm 1 accepts just all subsequences of T .

Proof: We prove the lemma in two steps.

1. If S is a subsequence of the string T then S is accepted by the automaton (induction by the length of S):

Step 1: $|S| = 1, S = s_1$. If s_1 occurs in T then state 0 of the automaton has transition labeled with s_1 and the automaton accepts S .

Step 2: A string $S_k = s_1 s_2 \dots s_k$ is a subsequence of T and is accepted by the automaton. Let us create a new string $S_{k+1} = s_1 s_2 \dots s_k s_{k+1}$ by adding a symbol s_{k+1} to the end of S_k . There exists a sequence i_1, i_2, \dots, i_k such that $s_1 = t_{i_1}, s_2 = t_{i_2}, \dots, s_k = t_{i_k}$ (the automaton will finish in state i_k after accepting S_k). If there exists i_{k+1} such that $i_k < i_{k+1} \leq n$ and $s_{k+1} = t_{i_{k+1}}$, then state i_k has transition labeled with s_{k+1} and the automaton accepts S_{k+1} .

2. If S is accepted by the automaton then S is a subsequence of T (induction by the length of S):

Step 1: $|S| = 1, S = s_1$. If S is accepted by the automaton then state 0 has the transition labeled with s_1 . State 0 has transition labeled with s_1 only if there exists $j, 1 \leq j \leq n$ such that $s_1 = t_j$.

Step 2: A string $S_k = s_1 s_2 \dots s_k$ is accepted by the automaton and there exists a sequence i_1, i_2, \dots, i_k such that $s_1 = t_{i_1}, s_2 = t_{i_2}, \dots, s_k = t_{i_k}$. We create a new string $S_{k+1} = s_1 s_2 \dots s_k s_{k+1}$ by adding a symbol s_{k+1} to the end of S_k . The automaton will finish in state i_k after accepting S_k . If state i_k has transition labeled s_{k+1} then there exists $i_{k+1}, i_k < i_{k+1} \leq n$ such that $s_{k+1} = t_{i_{k+1}}$. \triangle

Let us suppose an alphabet Σ and a text $T = t_1 t_2 \dots t_n$ over this alphabet. Each symbol a of Σ we encode using digits 0 and 1. For that we need at least $c = \lceil \log_2 z \rceil$ digits. The algorithm is incremental. When we add a new symbol encoded as $e_1 e_2 \dots e_c$, we need to ensure, that all previous final states have an output path labeled with $e_1 e_2 \dots e_c$. We use a binary tree as an auxiliary structure. The tree is built during the construction of the automaton. Each inner node of the tree has two lists (for 0 and for 1), which contents states, where ends the path labeled with the same symbols as the path in the tree, starting at any final state and is the longest possible. So, if a state s is in the list l_e in the node with the path $p_1 p_2 \dots p_q$ from the root, there exists a path in the DASG from any final state to state s , which is labeled $p_1 p_2 \dots p_q$ and state s has no output transition labeled e . We start with an automaton with the only state 0. States (tc) for $t = 0, 1, \dots, n$ are final. At the beginning, the tree has only the initial node.

```

1:  $l_0^\varepsilon \leftarrow \{0\}, l_1^\varepsilon \leftarrow \{0\}$ 
2: for  $k \leftarrow 0$  to  $(n - 1)$  do
3:   encode the symbol  $t_{(k+1)}$  as  $e_1 e_2 \dots e_c$ 
4:   set the root as the actual node in the tree
5:   for  $b \leftarrow 1$  to  $c$  do
6:     add state  $(ck + b)$ 
7:     for each state  $s$  in the list  $l_{e_b}$  in the actual node of the tree do
8:       add a transition labeled  $e_b$  between states  $s$  and  $(ck + b)$ 
9:       remove  $s$  from the list  $l_{e_b}$ 
10:    end for
11:    go to the child of the actual node of the tree through the transition
        labeled  $e_b$  and set the child as the actual node (if the child does not
        exist, create it and set both lists of new node empty)
12:    if  $b < c$  then
13:      add the state  $(ck + b)$  to the both lists in the actual node
14:    end if
15:  end for
16:  mark the state  $(ck + c)$  as a final state and add it to the both lists in
        the root of the tree
17: end for
    
```

Figure 6: Algorithm 2 (incremental construction of encoded DASG)

The algorithm is demonstrated in Fig. 7–11. The lists maintained in the node with the path p from the root are denoted as l_0^p and l_1^p , the symbols are coded this way: $a = 00$, $b = 01$, $c = 10$.

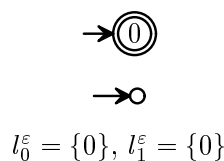


Figure 7: Encoded DASG for the string $T = \varepsilon$.

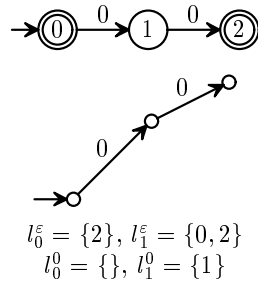


Figure 8: Encoded DASG for the string $T = a$.

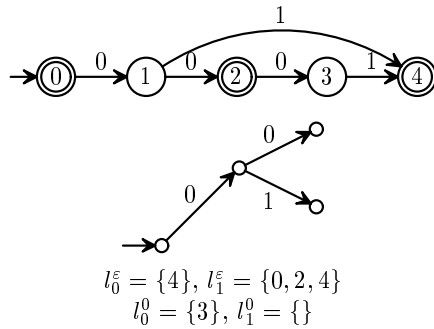


Figure 9: Encoded DASG for the string $T = ab$.

4.1 Number of transitions

The number of states grows to $1 + nc = 1 + n \lceil \log_2 n \rceil$. The maximum number of transitions is $c(2n - \frac{1}{2}(c + 1)) = \lceil \log_2 z \rceil (2n - \frac{1}{2}(\lceil \log_2 z \rceil + 1))$.

4.2 Complexity

The main cycle (line 2–15) is performed n times. Lines 1,3,4,6,11,12,13 and 16 require constant time. The cycle on line 5 is performed $\mathcal{O}(\log_2 z)$ times. The total number of transitions is $\mathcal{O}(n \log_2 z)$. Therefore, the total time complexity of lines 7–10 is $\mathcal{O}(n \log_2 z)$. Hence, the total time complexity is $\mathcal{O}(n \log_2 z)$.

The algorithm needs $\mathcal{O}(z + n \log_2 z)$ extra space for the tree and for the lists in its nodes. The subsequence test requires $\mathcal{O}(m \log_2 z)$ time in the worst case.

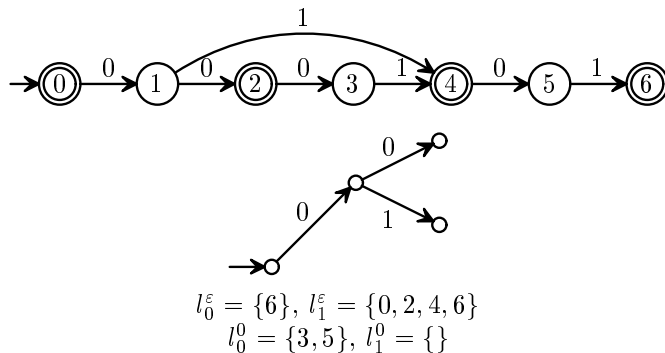
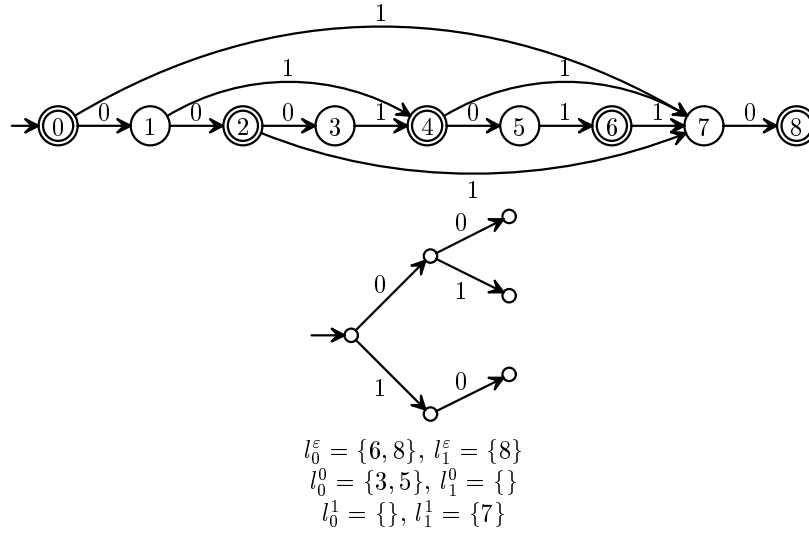


Figure 10: Encoded DASG for the string $T = abb$.


 Figure 11: Encoded DASG for the string $T = abc$.

5 DASG for a set of strings

Let us suppose an alphabet Σ and strings T_1, T_2, \dots, T_w over this alphabet. We extend Algorithm 1 to a set strings $\{T_1, T_2, \dots, T_w\}$. Let $L = \sum_{i=1}^w \text{length}(T_i)$.

The construction of DASG has two steps:

- Construction of inverted trie for reversed strings.
- Construction of an automaton.

Construction of inverted trie: Inverted trie arises from trie by reversing the transitions and can be constructed during the construction of trie (each node of trie will have one inverted transition). Final nodes of trie are initial nodes of inverted trie. Inverted trie is used as an auxiliary structure and served for finding common suffixes of the strings.

Construction of an automaton: For each string $T_i, 1 \leq i \leq w$ we will maintain:

- lists $l_a^i, a \in \Sigma$ of states, which have no output transition for the symbol a
- actual position act_i in inverted trie
- actual position $last_i$ in the automaton

In each node of inverted trie we save the number of corresponding state in the automaton. Let v denotes this number, let γ denotes the transition function in trie, and let δ denotes the transition function in the automaton. For each transition of the automaton we have to remember, which strings it belongs to. This set is denoted E . We start with the automaton with the only state 0. Each state of the automaton is final. *Set* is ordered set with two defined operation: *first* return the first string in the set and *next* return the successor of the string. The total number of states after each step is in the variable *states*.

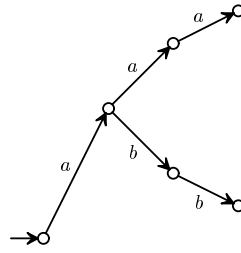
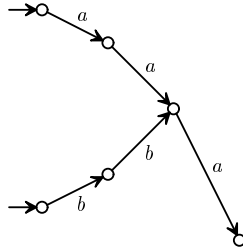
The algorithm is demonstrated in Fig. 15–18.

```

1: for  $i \leftarrow 1$  to  $w$  do
2:   for each  $a \in \Sigma$  do
3:      $l_a^i \leftarrow \{0\}$ 
4:   end for
5:    $act_i \leftarrow$  final state for the string  $T_i$  in trie
6:    $last_i \leftarrow 0$ 
7: end for
8: for each node  $i$  in trie do
9:    $v(i) \leftarrow 0$ 
10: end for
11:  $Set \leftarrow \{T_1, T_2, \dots, T_w\}$ 
12:  $states \leftarrow 1$ 
13:  $c \leftarrow 1$ 
14:  $p \leftarrow 1$ 
15: for  $k \leftarrow 1$  to  $L$  do
16:    $M \leftarrow \emptyset$ 
17:    $symbol \leftarrow$   $p$ -th symbol of  $T_c$ 
18:    $act_c \leftarrow \gamma(act_c, symbol)$ 
19:   if  $\delta(last_c, symbol) \neq \emptyset$  then
20:      $new\_state \leftarrow \delta(last_c, symbol)$ 
21:   else if  $v(act_c) > 0$  then
22:      $new\_state \leftarrow v(act_c)$ 
23:   else
24:     add state  $states$ 
25:      $new\_state \leftarrow states$ 
26:      $v(act_c) \leftarrow states$ 
27:      $states \leftarrow states + 1$ 
28:   end if
29:    $last_c \leftarrow new\_state$ 
30:   for each  $s \in l_{symbol}^c$  do
31:     if  $\delta(s, symbol) \neq \emptyset$  then
32:        $M \leftarrow M \cup \{\delta(s, symbol)\}$ 
33:        $E(s, symbol) \leftarrow E(s, symbol) \cup \{c\}$ 
34:     else
35:        $\delta(s, symbol) \leftarrow new\_state$ 
36:        $E(s, symbol) \leftarrow \{c\}$ 
37:     end if
38:   end for
39:    $l_{symbol}^c \leftarrow \emptyset$ 
40:    $M \leftarrow M \cup \{new\_state\}$ 
41:   for each  $a \in A$  do
42:      $l_a^c \leftarrow l_a^c \cup M$ 
43:   end for
44:   if  $p = length(T_c)$  then
45:      $Set \leftarrow Set \setminus \{T_c\}$ 
46:   end if
47:   if  $next(Set, c)$  is defined then
48:      $d$  is defined as follows:  $next(Set, c) = T_d$ 
49:   else
50:      $d$  is defined as follows:  $first(Set, c) = T_d$ 
51:      $p \leftarrow p + 1$ 
52:   end if
53:    $c \leftarrow d$ 
54: end for

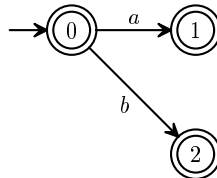
```

Figure 12: Algorithm 3 (extension of DASG for a set of strings $\{T_1, T_2, \dots, T_w\}$)


 Figure 13: Trie for the reversed strings aaa and bba .

 Figure 14: Inverted trie for the reversed strings aaa and bba .


$$l_a^1 = \{0\}, l_b^1 = \{0\}$$

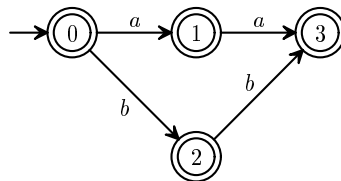
$$l_a^2 = \{0\}, l_b^2 = \{0\}$$

 Figure 15: Extension of DASG for the $Set = \{\varepsilon\}$.


$$l_a^1 = \{1\}, l_b^1 = \{0, 1\}$$

$$l_a^2 = \{0, 2\}, l_b^2 = \{2\}$$

$$E(0, a) = \{1\}, E(0, b) = \{2\}$$

 Figure 16: Extension of DASG for the $Set = \{a, b\}$.


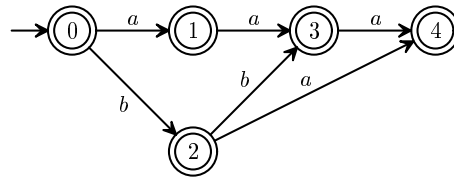
$$l_a^1 = \{3\}, l_b^1 = \{0, 1, 3\}$$

$$l_a^2 = \{0, 2, 3\}, l_b^2 = \{3\}$$

$$E(0, a) = \{1\}, E(0, b) = \{2\}$$

$$E(1, a) = \{3\}, E(2, b) = \{3\}$$

 Figure 17: Extension of DASG for the $Set = \{aa, bb\}$.



$$\begin{aligned}
 l_a^1 &= \{4\}, l_b^1 = \{0, 1, 3, 4\} \\
 l_a^2 &= \{1, 4\}, l_b^2 = \{1, 3, 4\} \\
 E(0, a) &= \{1\}, E(0, b) = \{2\} \\
 E(1, a) &= \{1\}, E(2, b) = \{2\} \\
 E(3, a) &= \{1\}, E(2, a) = \{2\}
 \end{aligned}$$

Figure 18: Extension of DASG for the $Set = \{aaa, bba\}$.

5.1 Number of states

For each symbol of the string, except for the last, a new state can be added. Hence, the DASG constructed in Algorithm 3 has at most $1 + \sum_{i=1}^w (\text{length}(T_i) - 1) + 1 = 2 + L - w$ states. DASG has this number of states if no two strings have any common nonempty prefix and suffix.

Each state can have at most z output transitions. Therefore, the total number of transitions is $\mathcal{O}(Lz)$.

5.2 Complexity

Construction of inverted trie requires $\mathcal{O}(L)$ time and $\mathcal{O}(L)$ extra space. For the time analysis of the Algorithm 3 is important the time complexity of set operations. We use four of them: *insert a member*, *delete a member*, *assign a value* and *union*. Suppose, that we use a balanced tree for the representation of the set M (another possibilities are a member function or a list). Then, the operations *assign* and *union* require $\mathcal{O}(|M|)$ time and the other operations require $\mathcal{O}(\log |M|)$ time.

Lines 1–7 require $\mathcal{O}(wz)$ time, lines 8–10 require $\mathcal{O}(L)$ time, line 11 requires $\mathcal{O}(w)$ time. The main cycle (lines 15–53) is performed L times. Line 16 requires $\mathcal{O}(L)$ time, and lines 18–20 require $\mathcal{O}(\log z)$ time. The cycle on the lines 30–38 is performed at most L times. Its time complexity is $\mathcal{O}(L \log L)$ (line 32 requires $\mathcal{O}(\log L)$ time). Line 40 requires $\mathcal{O}(\log L)$ time, lines 41–43 require $\mathcal{O}(Lz)$ time. Hence, the total time complexity is $\mathcal{O}(L^2 + L^2 \log L + L \log L + L^2 z) \approx \mathcal{O}(L^2 \log L)$ for arbitrary alphabet.

We need $\mathcal{O}(L)$ space for trie, $\mathcal{O}(L)$ space for the set M , and $\mathcal{O}(Lwz)$ space for the lists l_a^c . Hence, the total required extra space is $\mathcal{O}(Lwz)$.

5.3 Application: the longest common subsequence problem

The longest common subsequence (LCS) problem is known problem with applications in many areas. There are efficient algorithms that solve the LCS problem for two strings, for example [4].

Let us define the following problem (as in [2]): What is the longest common subsequence between any $k \leq w$ strings in a set S of w strings?

To solve this problem, we construct DASG for the set S and append to each transition $\delta(q, a)$ the number of strings in the set E (denoted as $num(q, a)$) and to

each state q the number of its input edges (denoted as c_q) and the number of its input edges with $num(q, a)$ greater or equal k (denoted as ck_q). We do not need the set E in this case. Then, we traverse DASG. During the traversing we use LIFO (Last-In-First-Out) memory as an auxiliary structure (denoted as $Stack$). Dot (\cdot) denotes concatenation. The longest sequence of input symbols from the initial state to the state q is stored in cs_q .

```

1:  $lcs \leftarrow \varepsilon$ 
2: for each state  $q$  do
3:    $cs_q \leftarrow \varepsilon$ 
4: end for
5:  $Stack \leftarrow 0$ 
6: while  $Stack$  is not empty do
7:    $q \leftarrow Pop$ 
8:   if  $length(cs_q) > length(lcs)$  then
9:      $lcs \leftarrow cs_q$ 
10:  end if
11:  for each  $a \in \Sigma$  such that  $\delta(q, a) \neq \emptyset$  do
12:     $r \leftarrow \delta(q, a)$ 
13:     $c_r \leftarrow c_r - 1$ 
14:    if  $c_r = 0$  then
15:       $Push(r)$ 
16:    end if
17:    if  $num(q, a) \geq k$  then
18:       $ck_r \leftarrow ck_r - 1$ 
19:      if  $ck_r = 0$  then
20:         $cs_r \leftarrow cs_q \cdot a$ 
21:      end if
22:    end if
23:  end for
24: end while

```

Figure 19: Algorithm 4 (the longest common subsequence)

The traversal of DASG requires $\mathcal{O}(Lz)$ time. For common subsequences cs_q we need $\mathcal{O}(Ly)$ space, where $y = \max\{length(T_i)\}$. Hence, the general longest common subsequence problem of w strings requires $\mathcal{O}(L^2 \log L + Lz)$ time for arbitrary alphabet. It is a better solution than presented in [2].

6 Conclusion

In section 3, we introduced a left-to-right algorithm for construction of DASG. It requires $\mathcal{O}(nz \log z)$ time and $\mathcal{O}(z)$ extra space for arbitrary alphabet. The subsequence test takes $\mathcal{O}(m \log z)$ time.

In section 4, we showed the modification of that algorithm for encoded DASG. The modified algorithm requires $\mathcal{O}(n \log z)$ time and $\mathcal{O}(z + n \log z)$ extra space. The subsequence test takes $\mathcal{O}(m \log z)$ time.

In section 5, we extended DASG for a set of strings and used it to solve the general longest common subsequence problem. Construction of DASG takes $\mathcal{O}(L^2 \log L)$ time and $\mathcal{O}(Lwz)$ extra space. The traversal of DASG requires $\mathcal{O}(Lz)$ time. Hence, the

solution of the general longest common subsequence problem requires $\mathcal{O}(L^2 \log L)$ total time and $\mathcal{O}(Lwz + Ly)$ space, where $y = \max\{\text{length}(T_i)\}$.

References

- [1] Baeza-Yates, R. A.: The Subsequence Graph of a Text. TAPSOFT'89, Proceedings of the International Joint Conference on Theory and Practice of Software Development, Volume 1: Advanced Seminar on Foundations of Innovative Software Development I and Colloquium on Trees in Algebra and Programming (CAAP'89), Lecture Notes in Computer Science 351, Barcelona, Spain, March 1989, pages 104–118.
- [2] Baeza-Yates, R. A.: Searching subsequences. Theoretical Computer Science 78 (1991), pages 363–378.
- [3] Crochemore, M., Rytter, W.: Text algorithms. Oxford University Press, 1994.
- [4] Hunt, J., Szymanski, T.: A fast algorithm for computing longest common subsequences. Communication of ACM 20, 1977, pages 350–353.
- [5] Hirschberg, D. S.: A linear space algorithm for computing maximal common subsequences. Communication of ACM, 18(6), 1975, pages 341–343.