



# Exact algorithms for the repetition-bounded longest common subsequence problem <sup>☆</sup>



Yuichi Asahiro <sup>a</sup>, Jesper Jansson <sup>b</sup>, Guohui Lin <sup>c</sup>, Eiji Miyano <sup>d,\*</sup>, Hirotaka Ono <sup>e</sup>,  
Tadatoshi Utashima <sup>d</sup>

<sup>a</sup> Kyushu Sangyo University, Fukuoka, Japan

<sup>b</sup> The Hong Kong Polytechnic University, Kowloon, Hong Kong

<sup>c</sup> University of Alberta, Edmonton, Canada

<sup>d</sup> Kyushu Institute of Technology, Iizuka, Japan

<sup>e</sup> Nagoya University, Nagoya, Japan

## ARTICLE INFO

### Article history:

Received 14 April 2020

Received in revised form 18 June 2020

Accepted 29 July 2020

Available online 4 August 2020

### Keywords:

Repetition-bounded longest common subsequence problem

Repetition-free

Exponential-time exact algorithms

Dynamic programming

NP-hardness

APX-hardness

## ABSTRACT

In this paper, we study exact, exponential-time algorithms for a variant of the classic LONGEST COMMON SUBSEQUENCE problem called the REPETITION-BOUNDED LONGEST COMMON SUBSEQUENCE problem (or RBLCS, for short): Let an *alphabet*  $S$  be a finite set of symbols and an *occurrence constraint*  $C_{occ}$  be a function  $C_{occ} : S \rightarrow \mathbb{N}$ , assigning an upper bound on the number of occurrences of each symbol in  $S$ . Given two sequences  $X$  and  $Y$  over the alphabet  $S$  and an occurrence constraint  $C_{occ}$ , the goal of RBLCS is to find a longest common subsequence of  $X$  and  $Y$  such that each symbol  $s \in S$  appears at most  $C_{occ}(s)$  times in the obtained subsequence. The special case where  $C_{occ}(s) = 1$  for every symbol  $s \in S$  is known as the REPETITION-FREE LONGEST COMMON SUBSEQUENCE problem (RFLCS) and has been studied previously; e.g., in [1], Adi et al. presented a simple (exponential-time) exact algorithm for RFLCS. However, they did not analyze its time complexity in detail, and to the best of our knowledge, there are no previous results on the running times of any exact algorithms for this problem. Without loss of generality, we will assume that  $|X| \leq |Y|$  and  $|X| = n$ . In this paper, we first propose a simpler algorithm for RFLCS based on the strategy used in [1] and show explicitly that its running time is  $O(1.44225^n)$ . Next, we provide a dynamic programming (DP) based algorithm for RBLCS and prove that its running time is  $O(1.44225^n)$  for any occurrence constraint  $C_{occ}$ , and even less in certain special cases. In particular, for RFLCS, our DP-based algorithm runs in  $O(1.41422^n)$  time, which is faster than the previous one. Furthermore, we prove NP-hardness and APX-hardness results for RBLCS on restricted instances.

© 2020 Elsevier B.V. All rights reserved.

<sup>☆</sup> A preliminary version of the paper appeared in Proceedings of the 13th Annual International Conference on Combinatorial Optimization and Applications (COCOA 2019), Vol. 11949 of Lecture Notes in Computer Science, Springer, pp. 1–12 (2019).

\* Corresponding author.

E-mail addresses: asahiro@is.kyusan-u.ac.jp (Y. Asahiro), jesper.jansson@polyu.edu.hk (J. Jansson), guohui@ualberta.ca (G. Lin), miyano@ces.kyutech.ac.jp (E. Miyano), ono@nagoya-u.jp (H. Ono), utashima.tadatoshi965@mail.kyutech.jp (T. Utashima).

$$\begin{aligned}
 X &= \text{TGACTCTGTGCA} \\
 Y &= \text{TGCTCAGTGCAC} \\
 Z &= \text{TGCTCGTA} \\
 Z' &= \text{TGAC}
 \end{aligned}$$

**Fig. 1.** If two sequences  $X$  and  $Y$ , and an occurrence constraint  $C_{occ}(A) = 1$ ,  $C_{occ}(C) = C_{occ}(G) = 2$  and  $C_{occ}(T) = 3$  are given as input, then  $Z$  is a repetition-bounded longest common subsequence satisfying the occurrence constraint  $C_{occ}$ . As another example, if  $C_{occ}(A) = C_{occ}(C) = C_{occ}(G) = C_{occ}(T) = 1$ , then  $Z'$  is a solution (i.e.,  $Z'$  is *repetition-free*).

## 1. Introduction

### 1.1. LONGEST COMMON SUBSEQUENCE problems

An *alphabet*  $S$  is a finite set of *symbols*. Let  $X$  be a sequence over the alphabet  $S$  and  $|X|$  be the length of the sequence  $X$ . For example,  $X = \langle x_1, x_2, \dots, x_n \rangle$  is a sequence of length  $n$ , where  $x_i \in S$  for  $1 \leq i \leq n$ , i.e.,  $|X| = n$ . For a sequence  $X = \langle x_1, x_2, \dots, x_n \rangle$ , another sequence  $Z = \langle z_1, z_2, \dots, z_c \rangle$  is a *subsequence* of  $X$  if there exists a strictly increasing sequence  $\langle i_1, i_2, \dots, i_c \rangle$  of indices of  $X$  such that for all  $j = 1, 2, \dots, c$ , we have  $x_{i_j} = z_j$ . Then, we say that a sequence  $Z$  is a *common subsequence* of  $X$  and  $Y$  if  $Z$  is a subsequence of both  $X$  and  $Y$ . Given two sequences  $X$  and  $Y$  as input, the goal of the LONGEST COMMON SUBSEQUENCE problem (LCS) is to find a *longest* common subsequence of  $X$  and  $Y$ .

LCS is a fundamental problem and has a long history [5,12,16,27]. The comparison of sequences via a longest common subsequence has been applied in several contexts where we want to find the maximum number of symbols that appear in the same order in two sequences. LCS is considered to be an important computational primitive in a variety of applications such as bioinformatics [3,4,20,22], data compression [26], spelling correction [21,27], and file comparison [2] since LCS plays a key role in measuring various types of sequence similarity.

LCS has been deeply investigated, and polynomial-time algorithms are well-known [16,17,22,23,27]. It is possible to generalize LCS to a set of three or more sequences; the goal is to compute a longest common subsequence of all input sequences. If the number of sequences is part of the input, then LCS of multiple sequences is NP-hard even on binary alphabet [19] and it is not approximable within factor  $O(n^{1-\varepsilon})$  on arbitrary alphabet for sequences of length  $n$  and any constant  $\varepsilon > 0$  [20]. Furthermore, some researchers introduced a constraint on the number of symbol occurrences in the solution. Bonizzoni et al. considered the EXEMPLAR LONGEST COMMON SUBSEQUENCE problem (ELCS) [10,24]. In ELCS, the alphabet  $S$  of symbols is divided into the mandatory alphabet  $S_m$  and the optional alphabet  $S_o$ , and ELCS restricts the numbers of symbol occurrences in  $S_m$  and  $S_o$  in the obtained solution. ELCS is APX-hard even for instances of two sequences [10]. In [11], Bonizzoni et al. proposed the following DOUBLY-CONSTRAINED LONGEST COMMON SUBSEQUENCE problem (DCLCS): Let a sequence constraint  $C$  be a set of sequences over an alphabet  $S$  and let an occurrence constraint  $C_{occ}$  be a function  $C_{occ} : S \rightarrow \mathbb{N}$ , assigning an upper bound on the number of occurrences of each symbol in  $S$ . Given two sequences  $X$  and  $Y$  over the alphabet  $S$ , a sequence constraint  $C$ , and an occurrence constraint  $C_{occ}$ , the goal of DCLCS is to find a longest common subsequence  $Z$  of  $X$  and  $Y$  such that each sequence in  $C$  is a subsequence of  $Z$  and  $Z$  contains at most  $C_{occ}(s)$  occurrences of each symbol  $s \in S$ . Bonizzoni et al. showed that DCLCS is NP-hard over an alphabet of three symbols [11].

Adi et al. introduced the REPETITION-FREE LONGEST COMMON SUBSEQUENCE problem (RFLCS) [1]: Given two sequences  $X$  and  $Y$  over an alphabet  $S$ , the goal of RFLCS is to find a “*repetition-free*” longest common subsequence of  $X$  and  $Y$ , where each symbol appears at most once in the obtained subsequence. In [1], Adi et al. proved that RFLCS is APX-hard even if each symbol appears at most twice in each of the given sequences.

### 1.2. Our new results

In this paper we study exact, exponential-time algorithms for RFLCS and its general form, called the REPETITION-BOUNDED LONGEST COMMON SUBSEQUENCE problem (RBLCS for short): Let  $S = \{s_1, s_2, \dots, s_k\}$  be an alphabet of  $k$  symbols. Recall that  $C_{occ}$  is an occurrence constraint  $C_{occ} : S \rightarrow \mathbb{N}$ , assigning an upper bound on the number of occurrences of each symbol in  $S$ . Given two sequences  $X$  and  $Y$  over the alphabet  $S$  and an occurrence constraint  $C_{occ}$ , the goal of RBLCS is to find a “*repetition-bounded*” longest common subsequence of  $X$  and  $Y$ , where each symbol  $s_i$  appears at most  $C_{occ}(s_i)$  times in the obtained subsequence for  $i = 1, 2, \dots, k$ . See Fig. 1 for examples. The special case where  $C_{occ}(s_i) = r$  for every  $i = 1, 2, \dots, k$  is referred to as the  $r$ -REPETITION-BOUNDED LONGEST COMMON SUBSEQUENCE problem ( $r$ -RBLCS for short). Note that the special case 1-RBLCS is identical to RFLCS. Also, it is easy to see that RBLCS is a special case of DCLCS where the sequence constraint  $C$  satisfies  $C = \emptyset$ . In [1], Adi et al. presented a simple (exponential-time) exact algorithm for 1-RBLCS, whose basic strategy is to enumerate all the subsequences consisting of representative symbols. However, they did not analyze its time complexity in detail as their focus was on establishing polynomial-time solvability and polynomial-time approximability. To the best of our knowledge, there are no previous results on the running times of any exact algorithms for this problem.

Without loss of generality, we will assume that  $|X| \leq |Y|$  and  $|X| = n$ . The contributions of this paper are summarized as follows:

1. We propose a simple algorithm for RFLCS based on the strategy used in [1] and show explicitly that its running time is  $O(1.44225^n)$ .
2. We provide a dynamic programming (DP) based algorithm for RBLCS and prove that its running time is  $O(1.44225^n)$  for any occurrence constraint  $C_{occ}$ , and even less in certain special cases. In particular, for RFLCS, our DP-based algorithm runs in  $O(1.41422^n)$  time, which is faster than the previous one.
3. The NP-hardness of RFLCS implies that RBLCS is NP-hard in general. In this paper we prove that for any integer  $r \geq 2$ ,  $r$ -RBLCS remains NP-hard even if each symbol appears exactly  $r$  or  $r + 1$  times in each of the given two sequences. Furthermore, we prove that  $r$ -RBLCS is APX-hard if every symbol appears exactly  $r$  or  $2r$  times in each of the given two sequences.

### 1.3. Related work

Although this paper focuses on the exact exponential algorithms, we here make a brief survey on previous results for RFLCS, from the viewpoints of heuristic, approximation and parameterized algorithms. In [1], Adi et al. introduced first heuristic algorithms for RFLCS. After that, several (meta)heuristic algorithms for RFLCS were proposed in [7,8,13,25]. A detailed comparison of those metaheuristic algorithms was given in [9]. As for the approximability of RFLCS, Adi et al. showed [1] that RFLCS admits an  $occ_{max}$ -approximation algorithm, where  $occ_{max}$  is the maximum number of occurrences of each symbol in one of the two input sequences. In [6], Blin et al. presented a randomized fixed-parameter algorithm for RFLCS parameterized by the size of the solution. In [15], Fernandes and Kiwi studied the asymptotic behavior of the length of a repetition-free longest common subsequence of two *random* sequences such that each symbol appears randomly, uniformly and independently.

### 1.4. Organization

The rest of the paper is organized as follows: Section 2 introduces notation which will be used throughout the paper, and then gives the formal definition of RBLCS. In Section 3, we present simple exact algorithms based on the strategy used in [1] for RFLCS and  $r$ -RBLCS and analyze their running times in detail. Then, we design the  $O(1.44225^n)$ -time DP-based algorithm for RBLCS in Section 4. Section 5 shows the NP-hardness and the APX-hardness of  $r$ -RBLCS on restricted instances. Finally, we conclude the paper in Section 6. The notation used throughout the paper is summarized in the appendix.

## 2. Preliminaries

Let  $S = \{s_1, s_2, \dots, s_k\}$  be an *alphabet*, i.e., a finite set of  $k$  symbols. Let  $X = \langle x_1, x_2, \dots, x_n \rangle$  and  $Y = \langle y_1, y_2, \dots, y_m \rangle$  be the given two sequences as input of RBLCS and  $Z$  be a common subsequence of  $X$  and  $Y$ . For the sequence  $X$ , the subsequence  $\langle x_i, \dots, x_j \rangle$  is denoted by  $X_{i..j}$ . Then, we define the  $i$ th *prefix* of  $X$ , for  $i = 1, \dots, n$ , as  $X_{1..i} = \langle x_1, x_2, \dots, x_i \rangle$ . Also, we define the  $i$ th *suffix* of  $X$ , for  $i = 1, \dots, n$ , as  $X_{i..n} = \langle x_i, x_{i+1}, \dots, x_n \rangle$ .  $X_{1..n}$  is  $X$ . Similarly, we define the  $j$ th prefix and the  $j$ th suffix of  $Y$ , for  $j = 1, \dots, m$ , as  $Y_{1..j} = \langle y_1, y_2, \dots, y_j \rangle$  and  $Y_{j..m} = \langle y_j, y_{j+1}, \dots, y_m \rangle$ , respectively.

Without loss of generality, we assume that both  $X$  and  $Y$  have all  $k$  symbols in  $S$ . Let  $occ(X, s_i)$ ,  $occ(Y, s_i)$  and  $occ(Z, s_i)$  be the numbers of occurrences of  $s_i$  in  $X$ ,  $Y$ , and  $Z$ , respectively, and thus  $occ(X, s_i) \geq 1$  and  $occ(Y, s_i) \geq 1$  for every symbol  $s_i$ . Let  $C_{occ}$  be an occurrence constraint. The REPETITION-BOUNDED LONGEST COMMON SUBSEQUENCE problem (RBLCS) can be formally defined as follows:

REPETITION-BOUNDED LONGEST COMMON SUBSEQUENCE problem (RBLCS)

**Input:** A pair of sequences  $X$  and  $Y$ , and an occurrence constraint  $C_{occ}$ .

**Goal:** Find a longest common subsequence  $Z$  of  $X$  and  $Y$  such that  $occ(Z, s) \leq C_{occ}(s)$  is satisfied for every  $s \in S$ .

We call such a sequence  $Z$  a *repetition-bounded* longest common subsequence. The special case where  $C_{occ}(s_i) = r$  for every  $i = 1, 2, \dots, k$  is referred to as the  $r$ -REPETITION-BOUNDED LONGEST COMMON SUBSEQUENCE problem ( $r$ -RBLCS). Also, 1-RBLCS is often called the REPETITION-FREE LONGEST COMMON SUBSEQUENCE problem (RFLCS).

When presenting the time complexity of algorithms, we often round the base of exponential functions up to the fifth digit after the decimal point. That is, for example, the running time  $O((\sqrt{2})^n)$  is written as  $O(1.41422^n)$  since  $\sqrt{2} = 1.414213562\dots$  and thus  $\sqrt{2} < 1.41422$ . Furthermore, since  $(\sqrt{2})^n poly(n)$  is sandwiched between  $1.41421^n$  and  $1.41422^n$  for every polynomial  $poly(n)$  of  $n$  and sufficiently large  $n$ , we write  $O((\sqrt{2})^n poly(n))$  as  $O(1.41422^n)$ .

## 3. Warm-up algorithms

In this section, we first focus on RFLCS, i.e., 1-RBLCS. The following brute-force exact algorithm for RFLCS obviously runs in  $O(2^n \cdot n \cdot m)$  time for two sequences  $X$  and  $Y$  where  $|X| = n$ ,  $|Y| = m$ , and  $|X| \leq |Y|$ : (i) First create all the subsequences of  $X$ , denoted by  $X_1$  through  $X_{2^n}$ . Then, (ii) obtain a longest common subsequence of  $X_i$  and  $Y$  for each  $i$

( $1 \leq i \leq 2^n$ ) by using an  $O(|X_i| \cdot m)$ -time algorithm for LCS [22,23,27]. Finally, (iii) find a repetition-free longest subsequence among those  $2^n$  common subsequences obtained in (ii) and output it.

In [1], Adi et al. presented the following algorithm for RFLCS, which is more sophisticated than the naive algorithm above: Let  $S$  be an alphabet of symbols. Suppose that each symbol in  $S_X \subseteq S$  appears in  $X$  fewer times than  $Y$ , and  $S_X = \{s_1, s_2, \dots, s_{|S_X|}\}$ . Also, let  $S_Y = S \setminus S_X$  and  $S_Y = \{s_{|S_X|+1}, s_{|S_X|+2}, \dots, s_{|S|}\}$ . (i) The algorithm creates all the subsequences  $X_1$  through  $X_{N_X}$  of the input sequence  $X$  such that all the symbols in  $S_X$  occur *exactly once*, but all the occurrences of symbols in  $S_Y$  are kept in  $X_i$  for every  $1 \leq i \leq N_X$ . Also, the algorithm creates all the subsequences  $Y_1$  through  $Y_{N_Y}$  of the input sequence  $Y$  such that all the symbols in  $S_Y$  occur *exactly once*, but all the occurrences of symbols in  $S_X$  are kept in  $Y_j$  for every  $1 \leq j \leq N_Y$ . Then, (ii) obtain a longest common subsequence of  $X_i$  and  $Y_j$  for every pair of  $i$  and  $j$  ( $1 \leq i \leq N_X$  and  $1 \leq j \leq N_Y$ ) by using an  $O(|X_i| \cdot |Y_j|)$ -time algorithm for the original LCS. Finally, (iii) find the longest subsequence among  $N_X \cdot N_Y$  common subsequences obtained in (ii), which must be repetition-free, and output it. Clearly, the running time of this method is  $O(N_X \cdot N_Y \cdot n \cdot m)$ . In [1], Adi et al. only claimed that if the number of symbols which appear twice or more in  $X$  and  $Y$  is bounded above by some constant, say,  $c$ , then the running time is  $O(m^c \cdot n \cdot m)$ , i.e., RFLCS is solvable in polynomial time. However, no upper bound on  $N_X \cdot N_Y$  was given in [1].

### 3.1. Repetition-free LCS

In this subsection we consider an algorithm called  $\text{ALG}$ , based on the same strategy as the one in [1] for RFLCS: (i) First create all the subsequences  $X_1$  through  $X_N$  of the input sequence  $X$  such that every symbol appears *exactly once* in  $X_i$  for  $1 \leq i \leq N$  in  $O(N \cdot n)$  time. Then, (ii) obtain a longest common subsequence of  $X_i$  and  $Y$  in  $O(n \cdot m)$  time for each  $i$  ( $1 \leq i \leq N$ ). Finally, (iii) find a repetition-free longest subsequence among  $N$  common subsequences obtained in (ii) and output it. Therefore, the running time of  $\text{ALG}$  is  $O(N \cdot n \cdot m)$ . It is important to note that  $\text{ALG}$  is identical to Adi et al.'s algorithm in [1] if  $S_X = S$  and thus  $S_Y = \emptyset$ .

A very simple argument gives us the first upper bound on  $N$  and the running time of  $\text{ALG}$ :

**Theorem 1.** *The running time of  $\text{ALG}$  is  $O(1.44467^n)$  for RFLCS on two sequences  $X$  and  $Y$  where  $|X| = n$ ,  $|Y| = m$ , and  $|X| \leq |Y|$ .*

**Proof.** Recall that  $X$  has  $k$  symbols,  $s_1$  through  $s_k$ , and  $s_i$  occurs  $\text{occ}(X, s_i)$  times in  $X$  for each integer  $i$ ,  $1 \leq i \leq k$ . Since the number  $N$  of subsequences in  $X$  created in (i) of  $\text{ALG}$  is bounded by the number of combinations of  $k$  symbols, the following is satisfied:

$$N \leq \prod_{i=1}^k \text{occ}(X, s_i). \tag{1}$$

From the inequality of arithmetic and geometric means, we have:

$$N \leq \left( \sum_{i=1}^k \text{occ}(X, s_i) / k \right)^k \leq (n/k)^k.$$

Here, by setting  $p \stackrel{\text{def}}{=} n/k \in \mathbb{R}^+$ , we have:

$$N \leq (p)^{n/p} = (p^{1/p})^n.$$

Note that the value of  $p^{1/p}$  becomes the maximum when  $p = e$ , where  $e$  denotes Euler's number. That is,  $N$  is bounded above by  $e^{n/e} < 1.444668^n$ . Therefore, the running time of  $\text{ALG}$  is  $O(1.444668^n \cdot n \cdot m) = O(1.44467^n)$ .  $\square$

A more refined estimate yields a smaller upper bound on  $N$ , which gives us the improved running time of  $\text{ALG}$ :

**Theorem 2.** *The running time of  $\text{ALG}$  is  $O(1.44225^n)$  for RFLCS on two sequences  $X$  and  $Y$  where  $|X| = n$ ,  $|Y| = m$ , and  $|X| \leq |Y|$ .*

**Proof.** Let  $\max_{1 \leq i \leq k} \{\text{occ}(X, s_i)\} = \text{occ}_{\max}$ . Also, let  $S_i = \{s_j \mid \text{occ}(X, s_j) = i\}$  for  $1 \leq i \leq \text{occ}_{\max}$ . That is,  $S_i$  is a set of symbols which appear exactly  $i$  times in  $X$ . Let  $n_i = i \times |S_i|$ . Since each symbol in  $S_i$  appears  $i$  times in  $X$ , the following equality holds:

$$\sum_{i=1}^{\text{occ}_{\max}} n_i = n. \tag{2}$$

In the following, we show a smaller upper bound on  $N$  than that in the proof of Theorem 1. From the fact that  $n_i = i \times |S_i|$ , one sees that the following equality holds:

**Table 1**  
Occurrence  $\text{occ}$  and running time  $T$ .

$\text{occ}$	2	3	4	5	6	7	8
$T$	$1.41422^n$	$1.44225^n$	$1.41422^n$	$1.37973^n$	$1.34801^n$	$1.32047^n$	$1.29684^n$

$$\prod_{i=1}^k \text{occ}(X, s_i) = \prod_{i=1}^{\text{occ}_{\max}} i^{n_i/i}. \tag{3}$$

Here, from the inequality of arithmetic and geometric means, the following is obtained:

$$\left( \prod_{i=1}^{\text{occ}_{\max}} (i^{1/i})^{n_i} \right)^{1/\sum_{i=1}^{\text{occ}_{\max}} n_i} \leq \frac{\sum_{i=1}^{\text{occ}_{\max}} (i^{1/i}) \cdot n_i}{\sum_{i=1}^{\text{occ}_{\max}} n_i}. \tag{4}$$

From the equations (1), (2), (3), and (4), we get:

$$N \leq \left( \frac{\sum_{i=1}^{\text{occ}_{\max}} (i^{1/i}) \cdot n_i}{n} \right)^n. \tag{5}$$

Now, it is important to note that  $i \in \mathbb{N}$ , i.e.,  $i$  is a positive integer while  $p = n/k$  is a positive real in the proof of the previous theorem. Therefore, by a simple calculation, one can verify that the following is true:

$$\max_{i \in \mathbb{N}} \{i^{1/i}\} = 3^{1/3}. \tag{6}$$

Hence, we can bound the number  $N$  of all the possible repetition-free common subsequences as follows:

$$\begin{aligned} N &\leq \left( \frac{\sum_{i=1}^{\text{occ}_{\max}} 3^{1/3} \cdot n_i}{n} \right)^n \\ &= \left( \frac{3^{1/3} \cdot \sum_{i=1}^{\text{occ}_{\max}} n_i}{n} \right)^n \\ &= (3^{1/3})^n \\ &< 1.4422496^n. \end{aligned}$$

As a result, the running time of our algorithm is  $O(1.4422496^n \cdot n \cdot m) = O(1.44225^n)$ . This completes the proof.  $\square$

**Corollary 1.** There is an  $O(\text{occ}^{n/\text{occ}} \cdot n \cdot m)$ -time algorithm to solve RFLCS for two sequences  $X$  and  $Y$  where  $|X| = n$ ,  $|Y| = m$ , and  $|X| \leq |Y|$  when every symbol occurs in  $X$  exactly  $\text{occ}$  times.

**Proof.** By the assumption,  $\text{occ} \times |S_{\text{occ}}| = n$  and thus  $|S_{\text{occ}'}| = 0$  for  $\text{occ} \neq \text{occ}'$ . From the inequality (5), one can easily obtain the following:

$$N \leq (\text{occ}^{1/\text{occ}})^n. \quad \square$$

Table 1 shows the running time  $T$  for each  $\text{occ} = 2, 3, \dots, 8$ .

### 3.2. $r$ -repetition-bounded LCS, $r \geq 2$

In this subsection we consider exact exponential algorithms for  $r$ -RBLCS. First, by a straightforward extension of the algorithm for RFLCS, we can design the following algorithm for  $r$ -RBLCS, named  $\text{ALG}_r$ : First, (i) create all the subsequences  $X_1$  through  $X_N$  of the input sequence  $X$  such that each symbol  $s$  appears *exactly*  $r$  times in  $X_i$  for  $1 \leq i \leq N$  if  $X$  has more than  $r$   $s$ 's; otherwise, all the occurrences of  $s$  in  $X$  are included in  $X_i$ . Each subsequence  $X_i$  can be created in  $O(n)$  time and thus the total running time of (i) is  $O(N \cdot n)$ . Then, (ii) obtain a longest common subsequence of  $X_i$  and  $Y$  in  $O(n \cdot m)$  time for each  $i$  ( $1 \leq i \leq N$ ). Finally, (iii) find a longest subsequence among  $N$  common subsequences obtained in (ii), which has at most  $r$  occurrences of every symbol, and output it. Therefore, the running time is  $O(N \cdot n \cdot m)$ .

Again, suppose that  $X$  has  $k$  symbols,  $s_1, s_2, \dots, s_k$ , and  $s_i$  occurs  $\text{occ}(X, s_i)$  times in  $X$  for each integer  $i$ ,  $1 \leq i \leq k$ , and  $\max_{1 \leq i \leq k} \{\text{occ}(X, s_i)\} = \text{occ}_{\max}$ . Let  $S_i = \{s_j \mid \text{occ}(X, s_j) = i\}$  for  $1 \leq i \leq \text{occ}_{\max}$  and  $n_i = i \times |S_i|$ . Then, we estimate an upper bound on  $N$  for each  $r$ :

**Table 2**  
 $N(r)$  and  $i$  for each  $r$ .

$r$	2	3	4	5	6	7	8	9	10
$N(r)$	1.58884	1.66852	1.72013	1.75684	1.78453	1.80630	1.82394	1.83856	1.85091
$i$	5	7	9	11	13	15	17	19	21

**Theorem 3.** For  $r$ -RBLCS on two sequences  $X$  and  $Y$  where  $|X| = n$ ,  $|Y| = m$ , and  $|X| \leq |Y|$ , the running time of  $ALG_r$  is as follows:

$$O \left( \left( \max_{i \in \mathbb{N}} \left\{ \left( \frac{i - \frac{r-1}{2}}{(r!)^{1/r}} \right)^{r/i} \right\} \right)^n \times n \cdot m \right).$$

**Proof.** First, the total number  $N$  of sequences created in (i) of  $ALG_r$  can be expressed as follows:

$$N = \prod_{i=1}^k \binom{OCC_i}{r} = \prod_{i=r+1}^{OCC_{max}} \binom{i}{r}^{n_i/i}.$$

From the inequality of arithmetic and geometric means, we can obtain the following inequality:

$$(i(i-1)(i-2)\dots(i-r+1))^{1/r} \leq \frac{(2i-r+1)r/2}{r} = i - \frac{r-1}{2}.$$

Therefore,  $N$  is bounded:

$$\begin{aligned} \prod_{i=r+1}^{OCC_{max}} \binom{i}{r}^{n_i/i} &\leq \prod_{i=r+1}^{OCC_{max}} \left( \frac{(i - \frac{r-1}{2})^r}{r!} \right)^{n_i/i} \\ &= \prod_{i=r+1}^{OCC_{max}} \left( \left( \frac{i - \frac{r-1}{2}}{(r!)^{1/r}} \right)^{r/i} \right)^{n_i} \\ &\leq \left( \max_{i \in \mathbb{N}} \left\{ \left( \frac{i - \frac{r-1}{2}}{(r!)^{1/r}} \right)^{r/i} \right\} \right)^n. \end{aligned}$$

This completes the proof.  $\square$

We have obtained the specific values of  $\max_{i \in \mathbb{N}} \left\{ \left( \frac{i - \frac{r-1}{2}}{(r!)^{1/r}} \right)^{r/i} \right\}$ , denoted by  $N(r)$ , and  $i$  for  $r$ -RBLCS by its empirical implementation. Table 2 shows  $N(r)$  and  $i$  for each  $r = 2, 3, \dots, 10$ .

#### 4. Dynamic programming algorithms for RBLCS

In this section we design a DP-based algorithm named DP for RBLCS.

##### 4.1. Original LCS

In this subsection, we briefly review the dynamic programming paradigm for the original LCS. For more details, e.g., see [14]. Let  $Z_{1..h} = \langle z_1, z_2, \dots, z_h \rangle$  be any longest common subsequence of the  $i$ th prefix  $X_{1..i}$  of  $X$  and the  $j$ th prefix  $Y_{1..j}$  of  $Y$ . It is well known that LCS has the following optimal-substructure property: (1) If  $x_i = y_j$ , then  $z_h = x_i = y_j$  and  $Z_{1..h-1}$  is a longest common subsequence of  $X_{1..i-1}$  and  $Y_{1..j-1}$ . (2) If  $x_i \neq y_j$ , then (a)  $z_h \neq x_i$  implies that  $Z_{1..h}$  is a longest common subsequence of  $X_{1..i-1}$  and  $Y_{1..j}$ ; (b)  $z_h \neq y_j$ , then  $Z_{1..h}$  is a longest common subsequence of  $X_{1..i}$  and  $Y_{1..j-1}$ .

We define  $L(i, j)$  to be the length of a longest common subsequence of  $X_{1..i}$  and  $Y_{1..j}$ . Then, the above optimal substructure of LCS gives the following recursive formula:

$$L(i, j) = \begin{cases} 0 & \text{if } i = 0 \text{ or } j = 0, \\ L(i-1, j-1) + 1 & \text{if } i, j > 0 \text{ and } x_i = y_j, \\ \max\{L(i, j-1), L(i-1, j)\} & \text{if } i, j > 0 \text{ and } x_i \neq y_j. \end{cases}$$

The DP algorithm for the original LCS computes each value of  $L(i, j)$  and stores it into a two-dimensional table  $L$  of size  $(n+1) \times (m+1)$  in row-major order.



In the case of RBLCS, we have to count the number of occurrences of every symbol in the prefix of  $Z$ . In the following we show a modified recursive formula and a DP-based algorithm for RBLCS.

#### 4.2. Repetition-bounded LCS

A trivial implementation of a dynamic programming approach might be to use the DP-based algorithm for LCS for multiple sequences: For RFLCS, we first generate all the permutations of  $k$  symbols, i.e.,  $k!$  repetition-free sequences of  $k$  symbols, say,  $X_1$  through  $X_{k!}$  and then obtain a longest common subsequence of  $X_i$ ,  $X$ , and  $Y$  for each  $i$  ( $1 \leq i \leq k!$ ) by using an  $O(|X_i| \cdot n \cdot m)$ -time DP-based algorithm solving LCS for multiple (three) sequences proposed in [18]. Therefore, the total running time is  $O(k! \cdot k \cdot n \cdot m)$ . For RBLCS, we first generate all the permutation of  $\sum_{i=1}^k C_{occ}(s_i)$  multiple symbols and then obtain a longest common subsequence  $Z$  such that  $occ(Z, s_i) \leq C_{occ}(s_i)$  is satisfied for every  $s_i \in S$ . Let  $N = \sum_{i=1}^k C_{occ}(s_i)$ . Then, the running time is  $O(N! \cdot N \cdot n \cdot m)$ , which is polynomial if  $N$  is constant.

In the following we design a faster DP-based algorithm DP. Let  $S_{>C_{occ}} = \{s_i \mid occ(X, s_i) > C_{occ}(s_i)\}$ . Now suppose that  $|S_{>C_{occ}}| = \ell$  and, without loss of generality,  $S_{>C_{occ}} = \{s_1, s_2, \dots, s_\ell\}$ . Then, we prepare an ‘‘occurrence’’ vector of length  $\ell$ , denoted by  $\mathbf{v} = (v_1, v_2, \dots, v_\ell)$ , where the  $p$ th component  $v_p$  corresponds to the  $p$ th symbol  $s_p$  for  $1 \leq p \leq \ell$  and  $v_p \in \{0, 1, \dots, C_{occ}(s_p)\}$ . Roughly speaking, DP uses the occurrence vector  $\mathbf{v}$  as an upper bound of occurrences of every symbol in an intermediate solution, in order not to break the occurrence constraint; it can therefore compute a subproblem of finding a repetition-bounded longest common subsequence of  $X_{1..i}$  and  $Y_{1..j}$ . Note that the number of possibilities in the occurrence vector is  $\prod_{i=1}^{\ell} (C_{occ}(s_i) + 1)$ .

For the occurrence vector  $\mathbf{v} = (v_1, v_2, \dots, v_{p-1}, v_p, v_{p+1}, \dots, v_\ell)$ , we define a new vector  $\mathbf{v}|_{p=q} = (v_1, v_2, \dots, v_{p-1}, q, v_{p+1}, \dots, v_\ell)$ . Note that if  $v_p = q$  in  $\mathbf{v}$ , then  $\mathbf{v}|_{p=q} = \mathbf{v}$ . Let  $\mathbf{0}$  be an  $\ell$ -dimensional 0-vector, i.e., the length of  $\mathbf{0}$  is  $\ell$  and all  $\ell$  components are 0. Also, let  $\mathbf{C}_{occ}$  be an  $\ell$ -dimensional vector such that the length of  $\mathbf{C}_{occ}$  is  $\ell$  and the  $p$ th component is  $C_{occ}(s_p)$  for  $1 \leq p \leq \ell$ .

Similarly to the previous subsection, we define  $L(i, j, \mathbf{v})$  to be the length of a repetition-bounded longest common subsequence of  $X_{1..i}$  and  $Y_{1..j}$  satisfying the occurrence vector  $\mathbf{v}$ , i.e., the length of the subsequence which does not break the occurrence constraint. Our algorithm for RBLCS computes each value of  $L(i, j, \mathbf{v})$  and stores it into a three-dimensional table  $L$  of size  $(n + 1) \times (m + 1) \times \prod_{i=1}^{\ell} (C_{occ}(s_i) + 1)$ .

**Theorem 4** (Optimal substructure of RBLCS). Consider the  $i$ th prefix  $X_{1..i}$  of  $X$  and the  $j$ th prefix  $Y_{1..j}$  of  $Y$ . Suppose that  $S_{>C_{occ}} = \{s_1, s_2, \dots, s_\ell\}$  be a set of  $\ell$  symbols such that each  $s_i$  occurs at least  $C_{occ}(s_i) + 1$  times in  $X$ . Let  $Z_{1..h} = \langle z_1, z_2, \dots, z_h \rangle$  be any repetition-bounded longest common subsequence of  $X_{1..i}$  and  $Y_{1..j}$  satisfying an occurrence vector  $\mathbf{v}$ . Then, the followings are satisfied:

- (1) If  $x_i = y_j = s_p$  and  $s_p \notin S_{>C_{occ}}$ , then  $z_h = s_p$  and  $Z_{1..h-1}$  is a repetition-bounded longest common subsequence of  $X_{1..i-1}$  and  $Y_{1..j-1}$  satisfying  $\mathbf{v}$ .
- (2) If  $x_i = y_j = s_p$ ,  $s_p \in S_{>C_{occ}}$  and  $v_p > 0$ , then  $z_h = s_p$  implies that  $Z_{1..h-1}$  is a repetition-bounded longest common subsequence of  $X_{1..i-1}$  and  $Y_{1..j-1}$  satisfying  $\mathbf{v}|_{p=v_p-1}$ .
- (3) If  $x_i = y_j = s_p$ ,  $s_p \in S_{>C_{occ}}$  and  $v_p = 0$ , then  $z_h \neq s_p$  and  $Z_{1..h}$  is a repetition-bounded longest common subsequence of  $X_{1..i-1}$  and  $Y_{1..j-1}$  satisfying  $\mathbf{v}$ .
- (4) If  $x_i \neq y_j$ , then
  - (a)  $z_h \neq x_i$  implies that  $Z_{1..h}$  is a repetition-bounded longest common subsequence of  $X_{1..i-1}$  and  $Y_{1..j}$  satisfying  $\mathbf{v}$ ;
  - (b)  $z_h \neq y_j$  implies that  $Z_{1..h}$  is a repetition-bounded longest common subsequence of  $X_{1..i}$  and  $Y_{1..j-1}$  satisfying  $\mathbf{v}$ .

**Proof.** We will verify (1) through (4):

(1) If  $z_h \neq x_i$ , then by appending  $x_i = y_j = s_p$  to  $Z_{1..h}$ , we can obtain a repetition-bounded common subsequence of  $X_{1..i}$  and  $Y_{1..j}$  of length  $h + 1$  satisfying  $\mathbf{v}$  since the number of  $s_p$ 's in  $Z$  is at most  $C_{occ}(s_p) - 1$  from the condition  $s_p \notin S_{>C_{occ}}$ . This contradicts the assumption that  $Z_{1..h}$  is a repetition-bounded longest common subsequence of  $X_{1..i}$  and  $Y_{1..j}$  satisfying  $\mathbf{v}$ . Therefore,  $z_h = x_i = y_j$  holds. What we have to do is to prove that the prefix  $Z_{1..h-1}$  is a repetition-bounded longest common subsequence of  $X_{1..i-1}$  and  $Y_{1..j-1}$  with length  $h - 1$  satisfying  $\mathbf{v}$ . For the purpose of obtaining a contradiction, suppose that there exists a repetition-bounded common subsequence  $Z'$  of  $X_{1..i-1}$  and  $Y_{1..j-1}$  with length greater  $h - 1$  satisfying  $\mathbf{v}$ . Then, by appending  $x_i = y_j = s_p$ , we obtain a repetition-bounded common subsequence of  $X_{1..i}$  and  $Y_{1..j}$  whose length is greater than  $h$  satisfying  $\mathbf{v}$ , which is a contradiction.

(2) If  $z_h = x_i = y_j = s_p$ , then  $Z_{1..h-1}$  is a repetition-bounded common subsequence of  $X_{1..i-1}$  and  $Y_{1..j-1}$  such that  $s_p$  appears at most  $v_p - 1$  times in  $Z_{1..h-1}$ . Suppose that there exists a repetition-bounded common subsequence  $Z'$  of  $X_{1..i-1}$  and  $Y_{1..j-1}$  with length greater than  $h - 1$  satisfying  $\mathbf{v}|_{p=v_p-1}$ . Since the  $p$ th component of  $\mathbf{v}|_{p=v_p-1}$  is  $v_p - 1$ , by appending  $x_i = y_j = s_p$  to  $Z'$ , we obtain a repetition-bounded common subsequence of  $X_{1..i}$  and  $Y_{1..j}$  whose length is greater than  $h$  satisfying  $\mathbf{v}$ , which contradicts the assumption that  $Z_{1..h}$  is a repetition-bounded longest common subsequence of  $X_{1..i}$  and  $Y_{1..j}$  satisfying  $\mathbf{v}$ .

(3) Suppose that there exists a repetition-bounded common subsequence  $Z'$  of  $X_{1..i-1}$  and  $Y_{1..j-1}$  with length greater than  $h$  satisfying  $\mathbf{v}$ . From  $z_h \neq s_p$ ,  $Z'$  is also a repetition-bounded common subsequence of  $X_{1..i}$  and  $Y_{1..j}$  satisfying  $\mathbf{v}$ ,

which again contradicts the assumption that  $Z_{1..h}$  is a repetition-bounded longest common subsequence of  $X_{1..i}$  and  $Y_{1..j}$  satisfying  $\mathbf{v}$ .

(4)(a) ((b), resp.) If there is a repetition-bounded common subsequence  $Z'$  of  $X_{1..i-1}$  and  $Y_{1..j}$  ( $X_{1..i}$  and  $Y_{1..j-1}$ , resp.) with length greater than  $h$  satisfying  $\mathbf{v}$ , then  $Z'$  would also be a repetition-bounded longest common subsequence of  $X_{1..i}$  and  $Y_{1..j}$  satisfying  $\mathbf{v}$ , contradicting the assumption that  $Z$  is a repetition-bounded longest common subsequence of  $X_{1..i}$  and  $Y_{1..j}$  satisfying  $\mathbf{v}$ .  $\square$

Then, we can obtain the following recursive formula:

$$L(i, j, \mathbf{v}) = \begin{cases} 0 & \text{if } i = 0 \text{ or } j = 0, \\ L(i - 1, j - 1, \mathbf{v}) + 1 & \text{if } i, j > 0, x_i = y_j = s_p, \text{ and } s_p \notin S_{>C_{occ}} \text{ (Case (1))}, \\ L(i - 1, j - 1, \mathbf{v}|_{p=v_p-1}) + 1 & \text{if } i, j > 0, x_i = y_j = s_p, s_p \in S_{>C_{occ}}, \text{ and } v_p > 0 \text{ (Case (2))}, \\ L(i - 1, j - 1, \mathbf{v}) & \text{if } i, j > 0, x_i = y_j = s_p, s_p \in S_{>C_{occ}}, \text{ and } v_p = 0 \text{ (Case (3))}, \\ \max\{L(i - 1, j, \mathbf{v}), L(i, j - 1, \mathbf{v})\} & \text{otherwise (Case (4)).} \end{cases}$$

Here is an outline of our algorithm DP, which computes each value of  $L(i, j, \mathbf{v})$  and stores it into a three-dimensional table  $L$  of size  $(n + 1) \times (m + 1) \times \prod_{i=1}^{\ell} (C_{occ}(s_i) + 1)$ : Initially, we set  $L(i, j, \mathbf{v}) = 0$  and  $pre(i, j, \mathbf{v}) = \text{null}$  for every  $i, j$ , and  $\mathbf{v}$ . Then, the algorithm DP fills entries from  $L(1, 1, \mathbf{0})$  to  $L(1, 1, \mathbf{C}_{occ})$ , then from  $L(1, 2, \mathbf{0})$  to  $L(1, 2, \mathbf{C}_{occ})$ , next from  $L(1, 3, \mathbf{0})$  to  $L(1, 3, \mathbf{C}_{occ})$ , etc. After filling all the entries in the first “two-dimensional plane”  $L(1, j, \mathbf{v})$ , the algorithm fills all the entries in the second two-dimensional plane  $L(2, j, \mathbf{v})$ , and so on. Finally, DP fills all the entries in the  $n$ th plane. The algorithm DP also maintains a three dimensional table  $pre$  of size  $(n + 1) \times (m + 1) \times \prod_{i=1}^{\ell} (C_{occ}(s_i) + 1)$  to help us construct an optimal repetition-bounded longest subsequence. The entry  $pre(i, j, \mathbf{v})$  points to the table entry corresponding to the optimal subproblem solution chosen when computing  $L(i, j, \mathbf{v})$ .

---

**Algorithm DP**

---

**Input:** Two sequences  $X = \langle x_1, x_2, \dots, x_n \rangle$  and  $Y = \langle y_1, y_2, \dots, y_m \rangle$ , and an occurrence constraint  $C_{occ}$ .

**Output:** Repetition-bounded longest common subsequence  $Z$  of  $X$  and  $Y$ .

**Initialization:** Find  $S_{>C_{occ}} = \{s_1, s_2, \dots, s_k\}$  and then set  $L(i, j, \mathbf{v}) = 0$  and  $pre(i, j, \mathbf{v}) = \text{null}$  for each  $i, j$ , and  $\mathbf{v}$ .

```

1. for i = 1 to n
2.   for j = 1 to m
3.     for v = 0 to Cocc
4.       /* Case (1) */
5.       if xi == yj == sp and sp ∉ S>Cocc
6.         L(i, j, v) = L(i - 1, j - 1, v) + 1
7.         pre(i, j, v) = "(i - 1, j - 1, v)"
8.       /* Case (2) */
9.       elseif xi == yj == sp, sp ∈ S>Cocc, and vp > 0
10.        L(i, j, v) = L(i - 1, j - 1, v|p=vp-1) + 1
11.        pre(i, j, v) = "(i - 1, j - 1, v|p=vp-1)"
12.      /* Case (3) */
13.      elseif xi == yj == sp, sp ∈ S>Cocc and vp = 0
14.        L(i, j, v) = L(i - 1, j - 1, v)
15.        pre(i, j, v) = "(i - 1, j - 1, v)"
16.      /* Case (4) (a) */
17.      elseif L(i - 1, j, v) ≥ L(i, j - 1, v)
18.        L(i, j, v) = L(i - 1, j, v)
19.        pre(i, j, v) = "(i - 1, j, v)"
20.      /* Case (4) (b) */
21.      else L(i, j, v) = L(i, j - 1, v)
22.        pre(i, j, v) = "(i, j - 1, v)"

```

**Termination:** Construct a repetition-bounded longest common subsequence  $Z$  based on two tables  $L$  and  $pre$ , and then output  $Z$ .

---

We bound the running time of DP:

**Theorem 5.** The running time of DP is  $O(\prod_{i=1}^{\ell} (C_{occ}(s_i) + 1) \cdot n \cdot m)$  for RBLCS on two sequences  $X$  and  $Y$  where  $|X| = n$ ,  $|Y| = m$ , and  $|X| \leq |Y|$ .

**Proof.** The algorithm DP for RBLCS computes each value of  $L(i, j, \mathbf{v})$  and stores it into the three-dimensional table  $L$  of size  $(n + 1) \times (m + 1) \times \prod_{i=1}^{\ell} (C_{occ}(s_i) + 1)$ . Clearly, each table entry takes  $O(1)$  time to compute. As a result, the running time of DP is  $O(\prod_{i=1}^{\ell} (C_{occ}(s_i) + 1) \cdot n \cdot m)$ .  $\square$



**Table 3**  
Running time  $T$  of DP for  $r$ -RBLCS.

$r$	2	3	4	5	6	7	8
$T$	$1.44225^n$	$1.41422^n$	$1.37973^n$	$1.34801^n$	$1.32047^n$	$1.29684^n$	$1.27652^n$

By showing that  $\prod_{i=1}^{\ell} (C_{occ}(s_i) + 1) \leq (3^{1/3})^n$  is satisfied, we obtain the following corollary:

**Corollary 2.** *The running time of DP is  $O(1.44225^n)$  for RBLCS on two sequences  $X$  and  $Y$  where  $|X| = n$ ,  $|Y| = m$ , and  $|X| \leq |Y|$ .*

**Proof.** Let  $|S_{>C_{occ}}| = \ell$  again. Recall that for every  $s_i \in S_{>C_{occ}}$ ,  $occ(X, s_i) > C_{occ}(s_i)$ . That is,  $occ(X, s_i) \geq C_{occ}(s_i) + 1$  since both  $occ(X, s_i)$  and  $C_{occ}(s_i)$  are integers. Therefore, the following is satisfied:

$$\sum_{i=1}^{\ell} (C_{occ}(s_i) + 1) \leq \sum_{i=1}^{\ell} occ(X, s_i) \leq n. \tag{7}$$

Let  $C_{max} = \max_{s_i \in S_{>C_{occ}}} \{C_{occ}(s_i)\}$  be the maximum of the occurrence constraint. Also, let  $u_i = |\{s_j \mid occ(X, s_j) = i\}|$  be the number of symbols which appear exactly  $i$  times in  $X$  for  $1 \leq i \leq C_{max}$ . One sees that the term  $\prod_{i=1}^{\ell} (C_{occ}(s_i) + 1)$  in the running time of DP can be rewritten as follows:

$$\prod_{i=1}^{\ell} (C_{occ}(s_i) + 1) = \prod_{i=2}^{C_{max}+1} ((i - 1) + 1)^{u_{i-1}} = \prod_{i=2}^{C_{max}+1} i^{u_{i-1}} = \prod_{i=2}^{C_{max}+1} i^{V_i}, \tag{8}$$

where the rightmost equality holds by setting  $V_i \stackrel{\text{def}}{=} i \times u_{i-1}$ . Then, we can show the following upper bound of the summation from  $V_2$  to  $V_{C_{max}+1}$  from the above inequality (7):

$$\sum_{i=2}^{C_{max}+1} V_i = \sum_{i=1}^{\ell} (C_{occ}(s_i) + 1) \leq n. \tag{9}$$

By combining the (in)equalities (6), (8), and (9), we can obtain the following upper bound on  $\prod_{i=1}^{\ell} (C_{occ}(s_i) + 1)$ :

$$\prod_{i=1}^{\ell} (C_{occ}(s_i) + 1) \leq \prod_{i=2}^{C_{max}+1} 3^{\frac{V_i}{3}} = (3^{1/3})^{\sum_{i=2}^{C_{max}+1} V_i} \leq (3^{1/3})^n.$$

Therefore, the running time of DP is  $O(1.44225^n)$  for RBLCS.  $\square$

The algorithm DP works a little bit faster for RFLCS:

**Corollary 3.** *The running time of DP is  $O(1.41422^n)$  for RFLCS on two sequences  $X$  and  $Y$  where  $|X| = n$ ,  $|Y| = m$ , and  $|X| \leq |Y|$ .*

**Proof.** It is enough to prepare the three-dimensional table  $L$  of size  $(n + 1) \times (m + 1) \times 2^\ell$  for RFLCS. Clearly, each table entry takes  $O(1)$  time to compute. As a result, the running time of DP is  $O(2^\ell \cdot n \cdot m)$ . Recall that the number  $|S_{>C_{occ}}|$  of symbols which appear at least twice in  $X$  is defined to be  $\ell$ . This implies that  $\ell \leq \frac{n}{2}$ . Therefore,  $2^\ell \leq 2^{n/2} < 1.414214^n$  is satisfied; the running time is  $O(1.41422^n)$  for RFLCS.  $\square$

The running time for  $r$ -RBLCS is as follows:

**Corollary 4.** *The running time of DP is  $O((r + 1)^{n/(r+1)} \cdot n \cdot m)$  for  $r$ -RBLCS on two sequences  $X$  and  $Y$  where  $|X| = n$ ,  $|Y| = m$ , and  $|X| \leq |Y|$ .*

**Proof.** We prepare a three-dimensional table  $L$  of size  $(n + 1) \times (m + 1) \times (r + 1)^\ell$  and each table entry takes  $O(1)$  time to compute. Clearly  $\ell \leq \frac{n}{r+1}$ , i.e.,  $(r + 1)^\ell \leq (r + 1)^{n/(r+1)}$  holds.  $\square$

Table 3 shows the running time  $T$  of DP for  $r$ -RBLCS,  $r = 2, 3, \dots, 8$ .

**5. Hardness of RBLCS**

The NP-hardness (or the APX-hardness) of RFLCS implies that RBLCS on general instances is also NP-hard. In this section, we investigate the computational complexity of RBLCS on restricted instances. First, we consider  $r$ -RBLCS where the instance is a pair of sequences  $X$  and  $Y$  such that each symbol appears exactly  $r$  or  $r + 1$  times for any integer  $r \geq 2$ . We can show the following hardness result:

**Theorem 6.** For any integer  $r \geq 2$ ,  $r$ -RBLCS is NP-hard even if  $occ(X, s_i) \in \{r, r + 1\}$  and  $occ(Y, s_i) \in \{r, r + 1\}$  hold for every symbol  $s_i \in S$ .

**Proof.** We prove that the NP-hardness of  $r$ -RBLCS by providing the polynomial-time reduction from RFLCS to  $r$ -RBLCS. Suppose that a pair of sequences  $X = \langle x_1, x_2, \dots, x_n \rangle$  and  $Y = \langle y_1, y_2, \dots, y_m \rangle$  is an instance of RFLCS such that every symbol appears at most twice in each of two sequences. Recall that RFLCS is NP-hard even if each symbol appears at most twice in each of the given two sequences [1]. Let  $S = \{s_1, s_2, \dots, s_k\}$ . Then, we construct the following pair of sequences  $X^r$  and  $Y^r$  as an instance of  $r$ -RBLCS:

$$X^r = \langle x_1, x_2, \dots, x_n, \overbrace{s_1, \dots, s_1}^{r-1}, \overbrace{s_2, \dots, s_2}^{r-1}, \dots, \overbrace{s_k, \dots, s_k}^{r-1} \rangle$$

$$Y^r = \langle y_1, y_2, \dots, y_m, \overbrace{s_1, \dots, s_1}^{r-1}, \overbrace{s_2, \dots, s_2}^{r-1}, \dots, \overbrace{s_l, \dots, s_l}^{r-1} \rangle$$

That is, the  $n$ th prefix of  $X^r$  ( $m$ th prefix of  $Y^r$ , resp.) is  $X$  ( $Y$ , resp.), the next  $r - 1$  symbols of  $X^r$  ( $Y^r$ , resp.) are  $r - 1$  duplicates of  $s_1$ , the next  $r - 1$  symbols of  $X^r$  ( $Y^r$ , resp.) are  $r - 1$  duplicates of  $s_2$ , etc. This completes the reduction, which can be clearly done in polynomial time. One sees that every symbol appears exactly  $r$  or  $r + 1$  times in each of the two sequences  $X^r$  and  $Y^r$ .

In the following, we show that there is a repetition-free common subsequence  $Z$  of  $X$  and  $Y$  of length at least  $c$  if and only if there is a common subsequence  $Z^r$  of  $X^r$  and  $Y^r$  such that the length  $|Z^r|$  is at least  $c + k(r - 1)$  under the constraint  $occ(Z^r, s_i) \leq r$  for every symbol  $s_i \in S$ .

(Only-if part) Suppose that  $Z = \langle z_1, z_2, \dots, z_c \rangle$  is an optimal solution for RFLCS when its instance is the pair of sequences  $X$  and  $Y$ . Clearly,

$$Z^r = \langle z_1, z_2, \dots, z_c, \overbrace{s_1, \dots, s_1}^{r-1}, \overbrace{s_2, \dots, s_2}^{r-1}, \dots, \overbrace{s_k, \dots, s_k}^{r-1} \rangle$$

is a common subsequence of  $X^r$  and  $Y^r$  such that  $occ(Z^r, s_i) \leq r$  since the  $c$ th prefix of  $Z^r$  is repetition-free. The length of  $Z^r$  is (at least)  $c + k(r - 1)$ .

(If part) Suppose that  $Z^*$  is a repetition-bounded longest common subsequence such that  $occ(Z^*, s_i) \leq r$  for every symbol  $s_i \in S$  and the length of  $Z^*$  is at least  $c + k(r - 1)$ . If the number of symbols whose  $r$  occurrences are included in  $Z^*$  is at most  $c - 1$ , then the length of  $Z^*$  must be less than  $c + k(r - 1)$  by the following calculation (since the remaining  $k - c + 1$  symbols appear at most  $r - 1$  times):

$$r(c - 1) + (r - 1)(k - c + 1) = c + k(r - 1) - 1 < c + k(r - 1).$$

That is, there are at least  $c$  symbols whose  $r$  occurrences are included in  $Z^*$ . Suppose that  $s_1^*$  through  $s_c^*$  appear  $r$  times in  $Z^*$ . Observe that each of the  $(n + 1)$ st suffix  $X^r_{n+1..n+k(r-1)}$  of  $X^r$  and the  $(m + 1)$ st suffix  $Y^r_{m+1..m+k(r-1)}$  has exactly  $r - 1$  occurrences of every symbol  $s_i \in S$ . This implies that the  $n$ th prefix  $X^r_{1..n}$  of  $X^r$  and the  $m$ th prefix  $Y^r_{1..m}$  of  $Y^r$  has one or two  $s_i^*$ 's for  $i = 1, 2, \dots, c$ . Suppose, for example, that  $c = 5$ , and  $X^r_{1..n}$  and  $Y^r_{1..m}$  have the following structure:

$$X^r_{1..n} = \langle \dots, s_1^*, \dots, s_2^*, \dots, s_1^*, \dots, s_3^*, s_4^*, \dots, s_2^*, \dots, s_5^* \rangle$$

$$Y^r_{1..m} = \langle \dots, s_1^*, \dots, s_2^*, s_1^*, \dots, s_3^*, \dots, s_4^*, \dots, s_2^*, s_5^*, \dots \rangle$$

Then, the seventh prefix of  $Z^*$  must be  $Z^*_{1..7} = \langle s_1^*, s_2^*, s_1^*, s_3^*, s_4^*, s_2^*, s_5^* \rangle$ . Here, one sees that (at least) the leftmost occurrence of every  $s_i^*$  for  $i = 1, 2, \dots, c$  must be included in both the  $n$ th prefix  $X^r_{1..n}$  of  $X^r$  and the  $m$ th prefix  $Y^r_{1..m}$  of  $Y^r$ . As a result, we can obtain a repetition-free common subsequence of  $X^r_{1..n} = X$  and  $Y^r_{1..m} = Y$  of length at least  $c$ . For the above example, a repetition-free subsequence  $\langle s_1^*, s_2^*, s_3^*, s_4^*, s_5^* \rangle$  of length  $c = 5$  can be obtained from  $X$  and  $Y$  by removing the second  $s_1^*$  and the second  $s_2^*$  from  $Z^*_{1..7}$ . This completes the proof.  $\square$

If every symbol appears more times, then we can prove the APX-hardness of  $r$ -RBLCS by providing a *gap-preserving reduction* from RFLCS to  $r$ -RBLCS:

**Theorem 7.** For a pair of sequences  $X$  and  $Y$  and any integer  $r \geq 2$ ,  $r$ -RBLCS is APX-hard even if  $occ(X, s_i) \in \{r, 2r\}$  and  $occ(Y, s_i) \in \{r, 2r\}$  hold for every symbol  $s_i \in S$ .

**Proof.** Again, suppose that a pair of sequences  $X = \langle x_1, x_2, \dots, x_n \rangle$  and  $Y = \langle y_1, y_2, \dots, y_m \rangle$  is an instance of RFLCS such that every symbol appears either once or twice in each of the two sequences. Then, we construct the following pair of sequences  $X^r$  and  $Y^r$  as an instance of  $r$ -RBLCS, which are different from the previous  $X^r$  and  $Y^r$  in the proof of Theorem 6:

$$X^r = \langle \overbrace{x_1, x_1, \dots, x_1}^r, \overbrace{x_2, x_2, \dots, x_2}^r, \dots, \overbrace{x_n, x_n, \dots, x_n}^r \rangle$$

$$Y^r = \langle \overbrace{y_1, y_1, \dots, y_1}^r, \overbrace{y_2, y_2, \dots, y_2}^r, \dots, \overbrace{y_n, y_n, \dots, y_n}^r \rangle$$

That is, the first  $r$  symbols in  $X^r$  ( $Y^r$ , resp.) are  $r$  duplicates of  $x_1$  ( $y_1$ , resp.), the next  $r$  symbols in  $X^r$  ( $Y^r$ , resp.) are  $r$  duplicates of  $x_2$  ( $y_2$ , resp.), etc. This completes the reduction, which can be clearly done in polynomial time. One sees that every symbol appears exactly  $r$  or  $2r$  times in each of the two sequences  $X^r$  and  $Y^r$ .

Let  $Z$  and  $Z^r$  be optimal solutions of RFLCS and  $r$ -RBLCS for the pairs  $(X, Y)$  and  $(X^r, Y^r)$ , respectively. Also, let  $\Gamma(n, m)$  be a parameter function of the instance pair  $(X, Y)$  such that  $\Gamma : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N}$ . Next, we show that the above reduction satisfies the two conditions of gap-preserving reductions: (i) If  $|Z| \geq \Gamma(n, m)$ , then  $|Z^r| \geq r \times \Gamma(n, m)$ , and (ii) if  $|Z| < (1 - \varepsilon)\Gamma(n, m)$  for a fixed small positive constant  $\varepsilon > 0$ , then  $|Z^r| < r \times (1 - \varepsilon)\Gamma(n, m)$ . In the following let  $Z = \langle z_1^*, z_2^*, \dots, z_c^* \rangle$  be an optimal solution for RFLCS when its instance is the pair of sequences  $X$  and  $Y$ .

(i) Suppose that  $|Z| = c$ , i.e.,  $c \geq \Gamma(n, m)$  holds. Then, we consider the following sequence  $Z^r$  when its instance is the pair of sequences  $X^r$  and  $Y^r$ :

$$Z^r = \langle \overbrace{z_1^*, z_1^*, \dots, z_1^*}^r, \overbrace{z_2^*, z_2^*, \dots, z_2^*}^r, \dots, \overbrace{z_c^*, z_c^*, \dots, z_c^*}^r \rangle$$

From the above reduction, it is clear that  $Z^r$  is a common subsequence of  $X^r$  and  $Y^r$  such that each symbol  $z_i^*$  appears at most  $r$  times and thus the length of  $Z^r$  is  $r \times c$ . Hence,  $|Z^r| \geq r \times c = r \times \Gamma(n, m)$  holds.

(ii) Suppose that the length  $|Z| = c$  of the optimal solution  $Z$  of RFLCS is less than  $(1 - \varepsilon)\Gamma(n, m)$ . Also, suppose for the purpose of obtaining a contradiction that  $Z^r$  consists of at least  $c + 1$  different symbols, say,  $z_1^*$  through  $z_{c+1}^*$ . For example, suppose that  $c = 8$  and  $Z^r$  consists of nine symbols  $z_1^*$  through  $z_9^*$  as follows:

$$Z^r = \langle z_1^*, z_2^*, z_3^*, z_1^*, z_4^*, z_2^*, z_5^*, z_4^*, z_6^*, z_7^*, z_8^*, z_7^*, z_9^* \rangle$$

We can assume that the leftmost occurrence of each  $z_i^*$  appears in the subscript order in  $Z^r$ , i.e., the first symbol is  $z_1^*$ , the second symbol is  $z_2^*$ , etc, as shown above. Then, one can verify that the above sequence  $Z^r$  includes  $\langle z_1^*, z_2^*, z_3^*, z_4^*, z_5^*, z_6^*, z_7^*, z_8^*, z_9^* \rangle$  as a repetition-free subsequence. More generally, the sequence  $Z'' = \langle z_1^*, z_2^*, \dots, z_{c+1}^* \rangle$  of length  $c + 1$  must be a repetition-free common subsequence of  $X$  and  $Y$ , which is a contradiction.

If the optimal solution  $Z^r$  of  $r$ -RBLCS has at most  $c$  different symbols, then the length of  $Z^r$  is at most  $r \times c$ , which is less than  $r \times (1 - \varepsilon)\Gamma(n, m)$ . This completes the proof.  $\square$

### 6. Conclusion

We have studied a new variant of the LONGEST COMMON SUBSEQUENCE problem, called the REPETITION-BOUNDED LONGEST COMMON SUBSEQUENCE problem (RBLCS), and its special problem, called the  $r$ -REPETITION-BOUNDED LONGEST COMMON SUBSEQUENCE problem ( $r$ -RBLCS). For  $r = 1$ , 1-RBLCS is known as the REPETITION-FREE LONGEST COMMON SUBSEQUENCE problem. We first showed that for 1-RBLCS there is a simple exact algorithm whose running time is  $O(1.44225^n)$ . Then, for RBLCS, we designed a DP-based exact algorithm whose running time is  $O(1.44225^n)$ . In particular, the DP-based algorithm can solve 1-RBLCS in  $O(1.41422^n)$  time. To see that reducing the time complexity from  $O(1.44225^n)$  to  $O(1.41422^n)$  can be of practical importance, consider for example the case of  $n = 100$  and observe that  $1.41422^{100}$  is seven times smaller than  $1.44225^{100}$ . Hence, a promising direction for future research is to design faster exact exponential-time algorithms for RBLCS. Another challenge is to develop efficient approximation algorithms for RBLCS.

### Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

### Acknowledgements

This work was partially supported by PolyU Fund 1-ZE8L, the Natural Sciences and Engineering Research Council of Canada, JST CREST JPMJR1402, and Grants-in-Aid for Scientific Research of Japan (KAKENHI) Grant Numbers JP17K00016, JP17K00024, JP17K19960 and JP17H01698.

## Appendix. Summary of notation

**Note:** Some symbols that appear only in a restricted context are not listed.

$S$	an alphabet, i.e., a finite set of symbols
$X, Y$	two input sequences of symbols
$Z$	a common subsequence of the given two sequences
$n$	the length $ X $ of $X$
$m$	the length $ Y $ of $Y$
$X_{1..i}$	the $i$ th prefix of a sequence $X$
$X_{j..n}$	the $j$ th suffix of a sequence $X$
$C$	a sequence constraint, i.e., a set of sequences over an alphabet $S$
$C_{occ}(s)$	an occurrence constraint, i.e., a function assigning an upper bound on the number of occurrences of each symbol $s \in S$
$occ(X, s)$	the number of occurrences of symbol $s \in S$ in $X$
$occ_{max}$	the maximum number of occurrences of all the symbols $s \in S$ in $X$ , i.e., $occ_{max} = \max_{s \in S} \{occ(X, s)\}$
$S_i$	a set of every symbol $s$ which appears exactly $i$ times in $X$ , i.e., $S_i = \{s \mid occ(X, s) = i\}$
$S_{>C_{occ}}$	a set of every symbol $s$ which appears more than $C_{occ}(s)$ times in $X$ , i.e., $S_{>C_{occ}} = \{s \mid occ(X, s) > C_{occ}(s)\}$
$e$	Euler's number
$\mathbb{R}^+$	a set of positive reals
$\mathbb{N}$	a set of positive integers

## References

- [1] S.S. Adi, M.D.V. Braga, C.G. Fernandes, C.E. Ferreira, F.V. Martinez, M.-F. Sagot, M.A. Stefanos, C. Tjandraatmadja, Y. Wakabayashi, Repetition-free longest common subsequence, *Discrete Appl. Math.* 158 (2010) 1315–1324.
- [2] A. Aho, J. Hopcroft, J. Ullman, *Data Structures and Algorithms*, Addison-Wesley, 1983.
- [3] S.F. Altschul, W. Gish, W. Miller, E.W. Myers, D.J. Lipman, Basic local alignment search tool, *J. Mol. Biol.* 215 (3) (1990) 403–410.
- [4] R. Beal, T. Afrin, A. Farheen, D. Adjeroh, A new algorithm for “the LCS problem” with application in compressing genome resequencing data, in: *Proc. BIBM*, 2015, pp. 69–74.
- [5] L. Bergroth, H. Hakonen, T. Raita, A survey of longest common subsequence algorithms, in: *Proc. SPIRE*, 2000, pp. 39–48.
- [6] G. Blin, P. Bonizzoni, R. Dondi, F. Sikora, On the parameterized complexity of the repetition free longest common subsequence problem, *Inf. Process. Lett.* 112 (7) (2012) 272–276.
- [7] C. Blum, M.J. Blesa, B. Calvo, Beam-ACO for the repetition-free longest common subsequence problem, in: *Proc. EA* 2013, 2014, pp. 79–90.
- [8] C. Blum, M.J. Blesa, Construct, merge, solve and adapt: application to the repetition-free longest common subsequence problem, in: *Proc. EvoCOP2016*, 2016, pp. 46–57.
- [9] C. Blum, M.J. Blesa, A comprehensive comparison of metaheuristics for the repetition-free longest common subsequence problem, *J. Heuristics* 24 (3) (2018) 551–579.
- [10] P. Bonizzoni, G. Della Vedova, R. Dondi, G. Fertin, R. Rizzi, S. Vialette, Exemplar longest common subsequence, *IEEE/ACM Trans. Comput. Biol. Bioinform.* 4 (4) (2007) 535–543.
- [11] P. Bonizzoni, G. Della Vedova, R. Dondi, Y. Pirola, Variants of constrained longest common subsequence, *Inf. Process. Lett.* 110 (20) (2010) 877–881.
- [12] L. Bulteau, F. Hüffner, C. Komusiewicz, R. Niedermeier, Multivariate algorithmics for NP-hard string problems: the algorithmics column by Gerhard J. Woeginger, *Bull. Eur. Assoc. Theor. Comput. Sci.* 114 (2014).
- [13] M. Castelli, S. Beretta, L. Vanneschi, A hybrid genetic algorithm for the repetition free longest common subsequence problem, *Oper. Res. Lett.* 41 (6) (2013) 644–649.
- [14] T.H. Cormen, C.E. Leiserson, R.L. Rivest, C. Stein, *Introduction to Algorithms*, 3rd ed., The MIT Press, 2009.
- [15] C.G. Fernandes, M. Kiwi, Repetition-free longest common subsequence of random sequences, *Discrete Appl. Math.* 210 (2016) 75–87.
- [16] D.S. Hirschberg, Algorithms for the longest common subsequence problem, *J. ACM* 24 (4) (1977) 664–675.
- [17] D.S. Hirschberg, A linear space algorithm for computing maximal common subsequences, *Commun. ACM* 18 (6) (1975) 341–343.
- [18] S.Y. Itoga, The string merging problem, *BIT* 21 (1) (1981) 20–30.
- [19] D. Maier, The complexity of some problems on subsequences and supersequences, *J. ACM* 25 (2) (1978) 322–336.
- [20] T. Jiang, M. Li, On the approximation of shortest common supersequences and longest common subsequences, *SIAM J. Comput.* 24 (5) (1995) 1122–1139.
- [21] H.L. Morgan, Spelling correction in systems programs, *Commun. ACM* 13 (2) (1970) 90–94.
- [22] S.B. Needleman, C.D. Wunsch, A general method applicable to the search for similarities in the amino acid sequence of two proteins, *J. Mol. Biol.* 48 (3) (1970) 443–453.
- [23] D. Sankoff, Matching sequences under deletion/insertion constraints, *Proc. Natl. Acad. Sci. USA* 69 (1) (1972) 4–6.
- [24] D. Sankoff, Genome rearrangement with gene families, *Bioinformatics* 15 (11) (1999) 909–917.
- [25] R.S. Mincu, A. Popa, Better heuristic algorithms for the repetition free LCS and other variants, in: *Proc. SPIRE*, 2018, pp. 297–310.
- [26] J.A. Storer, *Data Compression: Methods and Theory*, Computer Science Press, 1988.
- [27] R.A. Wagner, M.J. Fischer, The string-to-string correction problem, *J. ACM* 21 (1) (1974) 168–173.