

# Longest Common Extension with Recompression

Tomohiro I

Kyushu Institute of Technology, Japan

tomohiro@ai.kyutech.ac.jp

## Abstract

Given two positions  $i$  and  $j$  in a string  $T$  of length  $N$ , a *longest common extension (LCE) query* asks for the length of the longest common prefix between suffixes beginning at  $i$  and  $j$ . A compressed LCE data structure is a data structure that stores  $T$  in a compressed form while supporting fast LCE queries. In this article we show that the *recompression* technique is a powerful tool for compressed LCE data structures. We present a new compressed LCE data structure of size  $O(z \lg(N/z))$  that supports LCE queries in  $O(\lg N)$  time, where  $z$  is the size of Lempel-Ziv 77 factorization without self-reference of  $T$ . Given  $T$  as an uncompressed form, we show how to build our data structure in  $O(N)$  time and space. Given  $T$  as a grammar compressed form, i.e., an straight-line program of size  $n$  generating  $T$ , we show how to build our data structure in  $O(n \lg(N/n))$  time and  $O(n + z \lg(N/z))$  space. Our algorithms are deterministic and always return correct answers.

## 1 Introduction

Given two positions  $i$  and  $j$  in a text  $T$  of length  $N$ , a *longest common extension (LCE) query*  $\text{LCE}(i, j)$  asks for the length of the longest common prefix between suffixes beginning at  $i$  and  $j$ . Since LCE queries play a central role in many string processing algorithms (see text book [6] for example), efficient LCE data structures have been extensively studied. If we are allowed to use  $O(N)$  space, optimal  $O(1)$  query time can be achieved by, e.g., lowest common ancestor queries [1] on the suffix tree of  $T$ . However,  $O(N)$  space can be too expensive nowadays as the size of strings to be processed is quite large. Thus, recent studies focus on more space efficient solutions.

Roughly there are three scenarios: Several authors have studied tradeoffs among query time, construction time and data structure size [17, 5, 4, 19]; In [16], Prezza presented in-place LCE data structures showing that the memory space for storing  $T$  can be replaced with an LCE data structure while retaining optimal substring extraction time; LCE data structures working on grammar compressed representation of  $T$  were studied in [7, 2, 3, 15].

In this article we pursue the third scenario, which is advantageous when  $T$  is highly compressible. In grammar compression,  $T$  is represented by a Context Free Grammar (CFG) that generates  $T$  and only  $T$ . In particular CFGs in Chomsky normal form, called Straight Line Programs (SLPs), are often considered as any CFG can be easily transformed into an SLP without changing the order of grammar size. Let  $\mathcal{S}$  be an arbitrary SLP of size  $n$  generating  $T$ . Bille et al. [3] showed a Monte Carlo randomized data structure of  $O(n)$  space that supports LCE queries in  $O(\lg N + \lg^2 \ell)$  time, where  $\ell$  is the answer to the LCE query. Because their algorithm is based on Karp-Rabin fingerprints, the answer is correct w.h.p (with high probability). If we always expect correct answers, we have to verify fingerprints in preprocessing phase, spending either  $O(N \lg N)$  time (w.h.p.) and  $O(N)$  space or  $O(\frac{N^2}{n} \lg N)$  time (w.h.p.) and  $O(n)$  space.

For a deterministic solution, I et al. [7] proposed an  $O(n^2)$ -space data structure, which can be built in  $O(n^2 h)$  time and  $O(n^2)$  space from  $\mathcal{S}$ , and supports LCE queries in  $O(h \lg N)$  time, where  $h$  is the height of  $\mathcal{S}$ . As will be stated in Theorem 2, we outstrip this result.

Our work is most similar to that presented in [15]. They showed that the signature encoding [13] of  $T$ , a special kind of CFGs that can be stored in  $O(z \lg N \lg^* N)$  space, can support LCE queries in  $O(\lg N + \lg \ell \lg^* N)$  time, where  $z$  is the size of LZ77 factorization of  $T$ .<sup>1</sup> The signature encoding is based on locally consistent parsing technique, which determines the parsing of a string by local surrounding. A key property of the signature encoding is that any occurrence of the same substring of length  $\ell$  in  $T$  is guaranteed to be compressed in almost same way leaving only  $O(\lg \ell \lg^* N)$  discrepancies in its

<sup>1</sup>Note that there are several variants of LZ77 factorization. We refer to LZ77 factorization as the one that is known as the *f-factorization without self-reference* unless otherwise noted.

Input	Time	Space	Reference
$T$	$N f_{\mathcal{A}}$	$z \lg N \lg^* N$	Theorem 3 (1a) of [15]
$T$	$N$	$N$	Theorem 3 (1b) of [15]
$\mathcal{S}$	$n f_{\mathcal{A}} \lg N \lg^* N$	$n + z \lg N \lg^* N$	Theorem 3 (3a) of [15]
$\mathcal{S}$	$n \lg \lg n \lg N \lg^* N$	$n \lg^* N + z \lg N \lg^* N$	Theorem 3 (3b) of [15]
LZ77	$z f_{\mathcal{A}} \lg N \lg^* M$	$z \lg N \lg^* N$	Theorem 3 (2) of [15]
$T$	$N$	$N$	this work, Theorem 1
$\mathcal{S}$	$n \lg(N/n)$	$n + z \lg(N/z)$	this work, Theorem 2
LZ77	$z \lg^2(N/z)$	$z \lg(N/z)$	this work, Corollary 3

Table 1: Comparison of construction time and space between ours and [15], where  $N$  is the length of  $T$ ,  $\mathcal{S}$  is an SLP of size  $n$  generating  $T$ ,  $z$  is the size of LZ77 factorization of  $T$ , and  $f_{\mathcal{A}}$  is the time needed for predecessor queries on a set of  $z \lg N \lg^* N$  integers from an  $N$ -element universe.

surrounding. As a result, an LCE query can be answered by tracing the  $O(\lg \ell \lg^* N)$  surroundings created over two occurrences of the longest common extension. The algorithm is quite simple as we simply simulate the traversal of the derivation tree on the CFG while matching substrings by appearances of the common variables. Note that another cost  $O(\lg N)$  is needed to traverse the derivation tree of height  $O(\lg N)$  from the root.

In this article we show that CFGs created by the *recompression* technique exhibit a similar property that can be used to answer LCE queries in  $O(\lg N)$  time. In recent years recompression has been proved to be a powerful tool in problems related to grammar compression [8, 9] and word equations [10, 11]. The main component of recompression is to replace some pairs in a string with variables of the CFG. Although we use global information (like the frequencies of pairs in the string) to determine which pairs to be replaced, the pairing itself is done very locally, i.e., all occurrences of the pairs are replaced. Then we can show that any occurrence of the same substring in  $T$  is guaranteed to be compressed in almost same way leaving only  $O(\lg N)$  discrepancies in its surrounding. This leads to an  $O(\lg N)$ -time algorithm to answer LCE queries, improving the  $O(\lg N + \lg \ell \lg^* N)$ -time algorithm of [15]. We also improve the data structure size from  $O(z \lg N \lg^* N)$  to  $O(z \lg(N/z))$ .<sup>2</sup>

In [15], the authors proposed efficient algorithms to construct the signature encoding from various kinds of input as summarized in Table 1. We achieve better and cleaner complexities of construction from SLPs, which is important to plug them in compressed string processing in which we are to solve problems on SLPs without decompressing the string explicitly. However it should be noted that the data structures in [15] also support efficient text edit operations. We are not sure if our data structures can be efficiently dynamized.

Theorems 1 and 2 show our main results. Note that our data structure is a simple CFG of height  $O(\lg N)$  on which we can simulate the traversal of the derivation tree to a target position in constant time per move. Thus, it naturally supports  $\text{Extract}(i, \ell)$  queries, which asks for retrieving the substring  $T[i..i + \ell - 1]$ , in  $O(\lg N + \ell)$  time.

**Theorem 1.** *Given a string  $T$  of length  $N$ , we can compute in  $O(N)$  time and space a compressed representation of  $T$  of size  $O(z \lg(N/z))$  that supports  $\text{Extract}(i, \ell)$  in  $O(\lg N + \ell)$  time and LCE queries in  $O(\lg N)$  time.*

**Theorem 2.** *Given an SLP of size  $n$  generating a string  $T$  of length  $N$ , we can compute in  $O(n \lg(N/n))$  time and  $O(n + z \lg(N/z))$  space a compressed representation of  $T$  of size  $O(z \lg(N/z))$  that supports  $\text{Extract}(i, \ell)$  in  $O(\lg N + \ell)$  time and LCE queries in  $O(\lg N)$  time.*

Suppose that we are given the LZ77-compression of size  $z$  of  $T$  as an input. Since we can convert the input into an SLP of size  $O(z \lg(N/z))$  [18], we can apply Theorem 2 to the SLP and get the next corollary.

**Corollary 3.** *Given the LZ77-compression of size  $z$  of a string  $T$  of length  $N$ , we can compute in  $O(z \lg^2(N/z))$  time and  $O(z \lg(N/z))$  space a compressed representation of  $T$  of size  $O(z \lg(N/z))$  that supports  $\text{Extract}(i, \ell)$  in  $O(\lg N + \ell)$  time and LCE queries in  $O(\lg N)$  time.*

<sup>2</sup>We believe that the space complexities of [15] can be improved to  $O(z \lg(N/z) \lg^* N)$  by using the same trick we use in Lemma 14.

Technically, this work owes very much to two papers [9, 8]. For instance, our construction algorithm of Theorem 1 is essentially the same as the grammar compression algorithm based on recompression presented in [9]. Our contribution is in discovering the above mentioned property that can be used for fast LCE queries. Also, we use the property to upper bound the size of our data structure in terms of  $z$  rather than the smallest grammar size  $g^*$ . Since it is known that  $z \leq g^*$  holds, an upper bound in terms of  $z$  is preferable. Our construction algorithm of Theorem 2 owes to [8], in which the recompression technique solves the fully-compressed pattern matching problems. Basically our results can be obtained by applying the technique. However, we make some contributions on top of it: We give a new observation that simplifies the implementation and analysis of a component of recompression called **BComp** (see Section 4.1.2). Also, we show that we can improve the time complexity from  $O(n \lg N)$  to  $O(n \lg(N/n))$ .

## 2 Preliminaries

An alphabet  $\Sigma$  is a set of characters. A string over  $\Sigma$  is an element in  $\Sigma^*$ . For any string  $w \in \Sigma^*$ ,  $|w|$  denotes the length of  $w$ . Let  $\varepsilon$  be the empty string, i.e.,  $|\varepsilon| = 0$ . Let  $\Sigma^+ = \Sigma^* \setminus \{\varepsilon\}$ . For any  $1 \leq i \leq |w|$ ,  $w[i]$  denotes the  $i$ -th character of  $w$ . For any  $1 \leq i \leq j \leq |w|$ ,  $w[i..j]$  denotes the substring of  $w$  beginning at  $i$  and ending at  $j$ . For convenience, let  $w[i..j] = \varepsilon$  if  $i > j$ . For any  $0 \leq i \leq |w|$ ,  $w[1..i]$  (resp.  $w[|w| - i + 1..|w|]$ ) is called the prefix (resp. suffix) of  $w$  of length  $i$ . We say that a string  $x$  occurs at position  $i$  in  $w$  iff  $w[i..|x| - 1] = x$ . A substring  $w[i..j] = c^d$  ( $c \in \Sigma, d \geq 1$ ) of  $w$  is called a *block* iff it is a maximal run of a single character, i.e.,  $w[i - 1] \neq c$  and  $w[j + 1] \neq c$ .

The text on which LCE queries are performed is denoted by  $T \in \Sigma^*$  with  $N = |T|$  throughout this paper. We assume that  $\Sigma$  is an integer alphabet  $[1..N^{O(1)}]$  and the standard word RAM model with word size  $\Omega(\lg N)$ .

The size of our compressed LCE data structure is bounded by  $O(z \lg(N/z))$ , where  $z$  is the size of the LZ77 factorization of  $T$  defined as follows:

**Definition 4** (LZ77 factorization). *The factorization  $T = f_1 f_2 \cdots f_z$  is the LZ77 factorization of  $T$  iff the following condition holds: For any  $1 \leq i \leq z$ , let  $p_i = |f_1 f_2 \cdots f_{i-1}| + 1$ , then  $f_i = T[p_i]$  if  $T[p_i]$  does not appear in  $T[1..p_i - 1]$ , otherwise  $f_i$  is the longest prefix of  $T[p_i..N]$  that occurs in  $T[1..p_i - 1]$ .*

**Example 5.** *The LZ77 factorization of abaabaabb is a · b · a · aba · ab · b and  $z = 6$ .*

In this article, we deal with grammar compressed strings, in which a string is represented by a Context Free Grammar (CFG) generating the string only. In particular, we consider *Straight-Line Programs (SLPs)* that are CFGs in Chomsky normal form. Formally, an SLP that generates a string  $T$  is a triple  $\mathcal{S} = (\Sigma, \mathcal{V}, \mathcal{D})$ , where  $\Sigma$  is the set of characters (terminals),  $\mathcal{V}$  is the set of variables (non-terminals),  $\mathcal{D}$  is the set of deterministic production rules whose righthand sides are in  $\mathcal{V}^2 \cup \Sigma$ , and the last variable derives  $T$ .<sup>3</sup> Let  $n = |\mathcal{V}|$ . We treat variables as integers in  $[1..n]$  (which should be distinguishable from  $\Sigma$  by having extra one bit), and  $\mathcal{D}$  as an injective function that maps a variable to its righthand side. We assume that given any variable  $X$  we can access in  $O(1)$  time to the data space storing the information of  $X$ , e.g.,  $\mathcal{D}(X)$ . We refer to  $n$  as the size of  $\mathcal{S}$  since  $\mathcal{S}$  can be encoded in  $O(n)$  space. Note that  $N$  can be as large as  $2^{n-1}$ , and so, SLPs have a potential to achieve exponential compression.

We extend SLPs by allowing run-length encoded rules whose righthand sides are of the form  $X^d$  with  $X \in \mathcal{V}$  and  $d \geq 2$ , and call such CFGs *run-length SLPs (RLSLPs)*. Since a run-length encoded rule can be stored in  $O(1)$  space, we still define the size of an RLSLP by the number of variables.

Let us consider the derivation tree  $\mathcal{T}$  of an RLSLP  $\mathcal{S}$  that generates a string  $T$ , where we delete all the nodes labeled with terminals for simplicity. That is, every node in  $\mathcal{T}$  is labeled with a variable. The height of  $\mathcal{S}$  is the height of  $\mathcal{T}$ . We say that a sequence  $C = v_1 \cdots v_m$  of nodes is a *chain* iff the nodes are all adjacent in this order, i.e., the beginning position of  $v_{i+1}$  is the ending position of  $v_i$  plus one for any  $1 \leq i < m$ .  $C$  is labeled with the sequence of labels of  $v_1 \cdots v_m$ .

For any sequence  $p \in \mathcal{V}^*$  of variables, let  $val_{\mathcal{S}}(p)$  denote the string obtained by concatenating the strings derived from all variables in the sequence. We omit  $\mathcal{S}$  when it is clear from context. We say that  $p$  generates  $val(p)$ . Also, we say that  $p$  occurs at position  $i$  iff there is a chain that is labeled with  $p$  and begins at  $i$ .

The next lemma, which is somewhat standard for SLPs, also holds for RLSLPs.

**Lemma 6.** *For any RLSLP  $\mathcal{S}$  of height  $h$  generating  $T$ , by storing  $|val(X)|$  for every variable  $X$ , we can support  $\text{Extract}(i, \ell)$  in  $O(h + \ell)$  time.*

<sup>3</sup>We treat the last variable as the starting variable.

### 3 LCE data structure built from uncompressed texts

In this section, we prove Theorem 1. We basically show that the RLSP obtained by grammar compression algorithm based on recompression [8] can be used for fast LCE queries. In Subsection 3.1 we first review the recompression and introduce notation we use. In Subsection 3.2 we present a new characterization of recompression, which is a key to our contributions.

#### 3.1 TtoG: Grammar compression based on recompression

In [8] Jez proposed an algorithm **TtoG** to compute an RLSP of  $T$  in  $O(N)$  time based on the recompression technique.<sup>4</sup> Let  $\mathbf{TtoG}(T)$  denote the RLSP of  $T$  produced by **TtoG**. We use the term *letters* for variables introduced by **TtoG**. Also, we use  $c$  (rather than  $X$ ) to represent a letter.

**TtoG** consists of two different types of compression **BComp** and **PComp**, which stand for Block Compression and Pair Compression, respectively.

- **BComp**: Given a string  $w$  over  $\Sigma = [1..|w|]$ , **BComp** compresses  $w$  by replacing all blocks of length  $\geq 2$  with fresh letters. Note that **BComp** eliminates all blocks of length  $\geq 2$  in  $w$ . We can conduct **BComp** in  $O(|w|)$  time and space (see Lemma 7).
- **PComp**: Given an string  $w$  over  $\Sigma = [1..|w|]$  that contains no block of length  $\geq 2$ , **PComp** compresses  $w$  by replacing all pairs from  $\dot{\Sigma}\dot{\Sigma}$  with fresh letters, where  $(\dot{\Sigma}, \ddot{\Sigma})$  is a partition of  $\Sigma$ , i.e.,  $\Sigma = \dot{\Sigma} \cup \ddot{\Sigma}$  and  $\dot{\Sigma} \cap \ddot{\Sigma} = \emptyset$ . We can deterministically compute in  $O(|w|)$  time and space a partition of  $\Sigma$  by which at least  $(|w| - 1)/4$  pairs are replaced (see Lemma 8), and conduct **PComp** in  $O(|w|)$  time and space (see Lemma 9).

Let  $T_0$  be a sequence of letters obtained by replacing every character  $c$  of  $T$  with a letter generating  $c$ . Then **TtoG** compresses  $T_0$  by applying **BComp** and **PComp** by turns until the string gets shrunk into a single letter. Since **PComp** compresses a given string by a constant factor  $3/4$ , the height of  $\mathbf{TtoG}(T)$  is  $O(\lg N)$ , and the total running time can be bounded by  $O(N)$  (see Lemma 10).

In order to give a formal description we introduce some notation below. **TtoG** transforms level by level  $T_0$  into strings,  $T_1, T_2, \dots, T_{\hat{h}}$  where  $|T_{\hat{h}}| = 1$ . For any  $0 \leq h \leq \hat{h}$ , we say that  $h$  is the *level* of  $T_h$ . If  $h$  is even, the transformation from  $T_h$  to  $T_{h+1}$  is performed by **BComp**, and production rules of the form  $c \rightarrow \tilde{c}^d$  are introduced. If  $h$  is odd, the transformation from  $T_h$  to  $T_{h+1}$  is performed by **PComp**, and production rules of the form  $c \rightarrow \hat{c}\hat{c}$  are introduced. Let  $\Sigma_h$  be the set of letters appearing in  $T_h$ . For any even  $h$  ( $0 \leq h < \hat{h}$ ), let  $\ddot{\Sigma}_h$  denote the set of letters with which there is a block of length  $\geq 2$  in  $T_h$ . For any odd  $h$  ( $0 \leq h < \hat{h}$ ), let  $(\dot{\Sigma}_h, \ddot{\Sigma}_h)$  denote the partition of  $\Sigma_h$  used in **PComp** of level  $h$ .

The following four lemmas show how to conduct **BComp**, **PComp**, and therefore **TtoG**, efficiently, which are essentially the same as respectively Lemma 2, Lemma 5, Lemma 6, and Theorem 1, stated in [8]. We give the proofs for the sake of completeness.

**Lemma 7.** *Given a string  $w$  over  $\Sigma = [1..|w|]$ , we can conduct **BComp** in  $O(|w|)$  time and space.*

*Proof.* We first scan  $w$  in  $O(|w|)$  time and list all the blocks of length  $\geq 2$ . Each block  $c^d$  ( $c \in \Sigma, d \geq 2$ ) at position  $i$  is listed by a triple  $(c, d, i)$  of integers in  $\Sigma$ . Next we sort the list according to the pair of integers  $(c, d)$ , which can be done in  $O(|w|)$  time and space by radix sort. Finally, we replace each block  $c^d$  by a fresh letter based on the rank of  $(c, d)$ .  $\square$

For any string  $w \in \Sigma^*$  that contains no block of length  $\geq 2$ , let  $\mathbf{Freq}_w(c, \tilde{c}, 0)$  (resp.  $\mathbf{Freq}_w(c, \tilde{c}, 1)$ ) with  $c > \tilde{c} \in \Sigma$  denote the number of occurrences of  $c\tilde{c}$  (resp.  $\tilde{c}c$ ) in  $w$ . We refer to the list of non-zero  $\mathbf{Freq}_w(c, \tilde{c}, \cdot)$  sorted in increasing order of  $c$  as the *adjacency list* of  $w$ . Note that it is the representation of the weighted directed graph in which there are exactly  $\mathbf{Freq}_w(c, \tilde{c}, 0)$  (resp.  $\mathbf{Freq}_w(c, \tilde{c}, 1)$ ) edges from  $c$  to  $\tilde{c}$  (resp. from  $\tilde{c}$  to  $c$ ). Each occurrence of a pair in  $w$  is counted exactly once in the adjacency list. Then the problem of computing a good partition  $(\dot{\Sigma}, \ddot{\Sigma})$  of  $\Sigma$  reduces to maximum directed cut problem on the graph. Algorithm 1 is based on a simple greedy  $1/4$ -approximation algorithm of maximum directed cut problem.

**Lemma 8.** *Given the adjacency list of size  $m$  of a string  $w \in \Sigma^*$ , Algorithm 1 computes in  $O(m)$  time a partition  $(\dot{\Sigma}, \ddot{\Sigma})$  of  $\Sigma$  such that the number of occurrences of pairs from  $\dot{\Sigma}\dot{\Sigma}$  in  $w$  is at least  $(|w| - 1)/4$ .*

<sup>4</sup>Indeed, the paper shows how to compute an SLP of size  $O(g^* \lg(N/g^*))$ , where  $g^*$  is the smallest SLP size to generate  $T$ .

*Proof.* In the foreach loop, we first run a 1/2-approximation algorithm of maximum “undirected” cut problem on the adjacency list, i.e., we ignore the direction of the edges here. For each  $c$  in increasing order, we greedily determine whether  $c$  is added to  $\dot{\Sigma}$  or to  $\ddot{\Sigma}$  depending on  $\sum_{\tilde{c} \in \dot{\Sigma}} \text{Freq}(c, \tilde{c}, \cdot) \geq \sum_{\tilde{c} \in \ddot{\Sigma}} \text{Freq}(c, \tilde{c}, \cdot)$ . Note that  $\sum_{\tilde{c} \in \dot{\Sigma}} \text{Freq}(c, \tilde{c})$  (resp.  $\sum_{\tilde{c} \in \ddot{\Sigma}} \text{Freq}(c, \tilde{c})$ ) represents the number of edges between  $c$  and a character in  $\dot{\Sigma}$  (resp.  $\ddot{\Sigma}$ ). By greedy choice, at least half of the edges in subject become the ones connecting two characters each from  $\dot{\Sigma}$  and  $\ddot{\Sigma}$ . Hence, in the end,  $|E|$  becomes at least  $(|w| - 1)/2$ , where let  $E$  denote the set of edges between characters from  $\dot{\Sigma}$  and  $\ddot{\Sigma}$  (recalling that there are exactly  $|w| - 1$  edges). Since each edge in  $E$  corresponds to an occurrence of a pair from  $\dot{\Sigma}\ddot{\Sigma} \cup \ddot{\Sigma}\dot{\Sigma}$  in  $w$ , at least one of the two partitions  $(\dot{\Sigma}, \ddot{\Sigma})$  and  $(\ddot{\Sigma}, \dot{\Sigma})$  covers more than half of  $E$ . Hence we achieve our final bound  $|E|/2 = (|w| - 1)/4$  by choosing an appropriate partition at Line 7.

In order to see that Algorithm 1 runs in  $O(m)$  time, we only have to care about Line 3 and Line 7. We can compute  $\sum_{\tilde{c} \in \dot{\Sigma}} \text{Freq}(c, \tilde{c}, \cdot)$  and  $\sum_{\tilde{c} \in \ddot{\Sigma}} \text{Freq}(c, \tilde{c}, \cdot)$  by going through all  $\text{Freq}(c, \cdot, \cdot)$  for fixed  $c$  in the adjacency list, which are consecutive in the sorted list. Since each element of the list is used only once, the cost for Line 3 is  $O(m)$  in total. Similarly the computation at Line 7 can be done by going through the adjacency list again. Thus the algorithm runs in  $O(m)$  time.  $\square$

---

**Algorithm 1:** How to compute a partition of  $\Sigma$  for PComp to compress  $w$  by 3/4.

---

**Input:** Adjacency list of  $w \in \Sigma^*$ .

**Output:**  $(\dot{\Sigma}, \ddot{\Sigma})$  s.t. # occurrences of pairs from  $\dot{\Sigma}\ddot{\Sigma}$  in  $w$  is at least  $(|w| - 1)/4$ .

*/\* The information whether  $c \in \Sigma$  is in  $\dot{\Sigma}$  or  $\ddot{\Sigma}$  is written in the data space for  $c$ , which can be accessed in  $O(1)$  time. \*/*

```

1  $\dot{\Sigma} \leftarrow \ddot{\Sigma} \leftarrow \emptyset;$ 
2 foreach  $c \in \Sigma$  in increasing order do
3   if  $\sum_{\tilde{c} \in \dot{\Sigma}} \text{Freq}_w(c, \tilde{c}, \cdot) \geq \sum_{\tilde{c} \in \ddot{\Sigma}} \text{Freq}_w(c, \tilde{c}, \cdot)$  then
4      $\lfloor$  add  $c$  to  $\dot{\Sigma}$ ;
5   else
6      $\lfloor$  add  $c$  to  $\ddot{\Sigma}$ ;
7 if # occurrences of pairs from  $\dot{\Sigma}\ddot{\Sigma} <$  # occurrences of pairs from  $\ddot{\Sigma}\dot{\Sigma}$  then
8    $\lfloor$  switch  $\dot{\Sigma}$  and  $\ddot{\Sigma}$ ;
9 return  $(\dot{\Sigma}, \ddot{\Sigma});$ 

```

---

**Lemma 9.** *Given a string  $w$  over  $\Sigma = [1..|w|]$  that contains no block of length  $\geq 2$ , we can conduct PComp in  $O(|w|)$  time and space.*

*Proof.* We first compute the adjacency list of  $w$ . This can be easily done in  $O(|w|)$  time and space by sorting the  $|w| - 1$  size multiset  $\{(w[i], w[i + 1], 0) \mid 1 \leq i < |w|, w[i] > w[i + 1]\} \cup \{(w[i], w[i + 1], 1) \mid 1 \leq i < |w|, w[i] < w[i + 1]\}$  by radix sort. Then by Lemma 8 we compute a partition  $(\dot{\Sigma}, \ddot{\Sigma})$  in linear time in the size of the adjacency list, which is  $O(|w|)$ . Next we scan  $w$  in  $O(|w|)$  time and list all the occurrences of pairs to be compressed. Each pair  $\tilde{c}\hat{c} \in \dot{\Sigma}\ddot{\Sigma}$  at position  $i$  is listed by a triple  $(\acute{c}, \grave{c}, i)$  of integers in  $\Sigma$ . Then we sort the list according to the pair of integers  $(\acute{c}, \grave{c})$ , which can be done in  $O(|w|)$  time and space by radix sort. Finally, we replace each pair with a fresh letter based on the rank of  $(\acute{c}, \grave{c})$ .  $\square$

**Lemma 10.** *Given a string  $T$  over  $\Sigma = [1..N^{O(1)}]$ , we can compute TtoG( $T$ ) in  $O(N)$  time and space.*

*Proof.* We first compute  $T_0$  in  $O(N)$  by sorting the characters used in  $T$  and replacing them with ranks of characters. Then we compress  $T_0$  by applying BComp and PComp by turns and get  $T_1, T_2 \dots T_{\hat{h}}$ . One technical problem is that characters used in an input string  $w$  of BComp and PComp should be in  $[1..|w|]$ , which is crucial to conduct radix sort efficiently in  $O(|w|)$  time (see Lemmas 7 and 9). However letters in  $T_h$  do not necessarily hold this property. To overcome this problem, during computation we maintain ranks of letters among those used in the current  $T_h$ , which should be in  $[1..|T_h|]$ , and use the ranks instead of letters for radix sort. If we have such ranks in each level, we can easily maintain them by radix sort for the next level. Now, in every level  $h$  ( $0 \leq h < \hat{h}$ ) the compression from  $T_h$  to  $T_{h+1}$  can be conducted in  $O(|T_h|)$  time and space. Since PComp compresses a given string by a constant factor, the total running time can be bounded by  $O(N)$  time.  $\square$

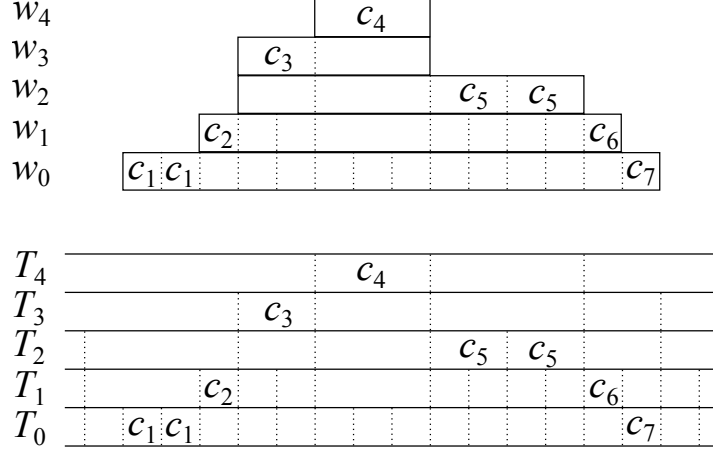


Figure 1: Suppose that  $c_1, c_7 \in \ddot{\Sigma}_0, c_2 \in \dot{\Sigma}_1, c_6 \in \dot{\Sigma}_1, c_5 \in \ddot{\Sigma}_2, c_3 \in \dot{\Sigma}_3$  and  $c_4 \in \ddot{\Sigma}_4$ . Then  $PSeq(w) = c_1c_1c_2c_3c_4c_5c_5c_6c_7$ .  $PSeq(w)$  occurs everywhere  $w$  occurs.

### 3.2 Popped sequences

We give a new characterization of recompression, which is a key to fast LCE queries as well as the upper bound  $O(z \lg(N/z))$  for the size of  $\text{TtoG}(T)$ . For any substring  $w$  of  $T$ , we define the *Popped Sequence* ( $PSeq$ ) of  $w$ , denoted by  $PSeq(w)$ .  $PSeq(w)$  is a sequence of letters such that  $val(PSeq(w)) = w$  and consists of  $O(\lg N)$  blocks of letters. It is not surprising that any substring can be represented by  $O(\lg N)$  blocks of letters because the height of  $\text{TtoG}(T)$  is  $O(\lg N)$ . The significant property of  $PSeq(w)$  is that it occurs at “every” occurrence of  $w$  (even if the occurrences overlap). A similar property has been observed in CFGs produced by locally consistent parsing and utilized for compressed indexes [12, 14] and a dynamic compressed LCE data structure [15]. For example, in [14, 15] the sequence having such a property is called the *common sequence* of  $w$  but its representation size is  $O(\lg |w| \lg^* N)$  rather than  $O(\lg N)$ .

$PSeq(w)$  is the sequence of letters characterized by the following procedure. Let  $w_0$  be the substring of  $T_0$  that generates  $w$ . We consider applying  $\text{BComp}$  and  $\text{PComp}$  to  $w_0$  exactly as we did to  $T$  but in each level we *pop* some letters out if the letters can be coupled with letters outside the scope. Formally, in increasing order of  $h \geq 0$ , we get  $w_{h+1}$  from  $w_h$  as follows:

- If  $h$  is even. We first pop out the leftmost and rightmost blocks of  $w_h$  if they are blocks of letter  $c \in \ddot{\Sigma}_h$ . Then we get  $w_{h+1}$  by applying  $\text{BComp}$  to the remaining string.
- If  $h$  is odd. We first pop out the leftmost letter and rightmost letter of  $w_h$  if they are letters in  $\dot{\Sigma}_h$  and  $\dot{\Sigma}_h$ , respectively. Then we get  $w_{h+1}$  by applying  $\text{PComp}$  to the remaining string.

We iterate this until the string disappears.  $PSeq(w)$  is the sequence obtained by concatenating the popped-out letters/blocks in an appropriate order. Note that for any occurrence of  $w$  the letters inside the  $PSeq(w)$  are compressed in the same way. Hence  $w_h$  is created for every occurrence of  $w$  and the occurrence of  $PSeq(w)$  is guaranteed (see also Figure 1).

The next lemma formalizes the above discussion.

**Lemma 11.** *For any substring  $w$  of  $T$ ,  $PSeq(w)$  consists of  $O(\lg N)$  blocks of letters. In addition,  $w$  occurs at position  $i$  iff  $PSeq(w)$  occurs at  $i$ .*

**Lemma 12.** *For any chain  $C$  whose label consists of  $m$  blocks of letters, the number of ancestor nodes of  $C$  is  $O(m)$ .*

*Proof.* Since a block is compressed into one letter, the number of parent nodes of  $C$  is at most  $m$ . As every internal node has two or more children, it is easy to see that there are  $O(m)$  ancestor nodes of the parent nodes of  $C$ .  $\square$

**Corollary 13.** *For any chain  $C$  corresponding to  $PSeq(T[b..e])$  for some interval  $[b..e]$ , the number of ancestor nodes of  $C$  is  $O(\lg N)$ .*

**Lemma 14.** *The size of  $\text{TtoG}(T)$  is  $O(z \lg(N/z))$ , where  $z$  is the size of the LZ77 factorization of  $T$ .*

*Proof.* We first show the bound  $O(z \lg N)$  and improve the analysis to  $O(z \lg(N/z))$  later.

Let  $f_1 \dots f_z$  be the LZ77 factorization of  $T$ . For any  $1 \leq i \leq z$ , let  $L_i$  be the set of letters used in the ancestor nodes of leaves corresponding to  $f_1 f_2 \dots f_i$ . Clearly  $|L_1| = O(\lg N)$ . For any  $1 < i \leq z$ , we estimate  $|L_i \setminus L_{i-1}|$ . Since  $f_i$  occurs in  $f_1 \dots f_{i-1}$ , we can see that the letters of  $PSeq(f_i)$  are in  $L_{i-1}$  thanks to Lemma 11. Let  $C_i$  be the chain corresponding to the occurrence  $|f_1 \dots f_{i-1} + 1|$  of  $PSeq(f_i)$ . Then, the letters in  $L_i \setminus L_{i-1}$  are only in the labels of ancestor nodes of  $C_i$ . Since  $PSeq(f_i)$  consists of  $O(\lg N)$  blocks of letters,  $|L_i \setminus L_{i-1}|$  is bounded by  $O(\lg N)$  due to Lemma 12. Therefore the size of  $\text{TtoG}(T)$  is  $O(z \lg N)$ .

In order to improve the bound to  $O(z \lg(N/z))$ , we employ the same trick that has been used in the literature. Let  $h = 2 \lg_{4/3}(N/z) = 2 \lg_{3/4}(z/N)$ . Recall that  $\text{PComp}$  compresses a given string by a constant factor  $3/4$ . Since  $\text{PComp}$  has been applied  $h/2$  times until the level  $h$ ,  $|T_h| \leq N(3/4)^{h/2} = z$ , and hence, the number of letters introduced in level  $\geq h$  is bounded by  $O(z)$ . Then, we can ignore all the letters introduced in level  $\geq h$  in the analysis of the previous paragraph, and by doing so, the bound  $O(\lg N)$  of  $|L_i \setminus L_{i-1}|$  is improved to  $O(h) = O(\lg(N/z))$ . This yields the bound  $O(z \lg(N/z))$  for the size of  $\text{TtoG}(T)$ .  $\square$

**Lemma 15.** *Given  $\text{TtoG}(T)$ , we can answer  $\text{LCE}(i, j)$  in  $O(\lg N)$  time.*

*Proof.* Let  $w$  be the longest common prefix of two suffixes beginning at  $i$  and  $j$ . In the light of Lemma 11,  $PSeq(w)$ , which consists of  $O(\lg N)$  blocks of letters, occurs at both  $i$  and  $j$ . Let  $C_i$  (resp.  $C_j$ ) be the chain that is labeled with  $PSeq(w)$  and begins at  $i$  (resp.  $j$ ). We can compute  $|w|$  by traversing the ancestor nodes of  $C_i$  and  $C_j$  simultaneously and matching  $PSeq(w)$  written in the labels of  $C_i$  and  $C_j$ . Note that we do matching from left to right and we do not have to know  $|w|$  in advance. Also, matching a block of letters in  $PSeq(w)$  can be done in  $O(1)$  time on run-length encoded rules. By Corollary 13, the number of ancestor nodes we have to visit is bounded by  $O(\lg N)$ . Thus, we get the lemma.  $\square$

### 3.3 Proof of Theorem 1

*Proof of Theorem 1.* By Lemma 10 we can compute  $\text{TtoG}(T)$  in  $O(N)$  time and space. Since the height of  $\text{TtoG}(T)$  is  $O(\lg N)$ , we can support  $\text{Extract}$  queries in  $O(\lg N + L)$  time due to Lemma 6.  $\text{LCE}$  queries can be supported in  $O(\lg N)$  time by Lemma 15.  $\square$

## 4 LCE data structure built from SLPs

In this section, we prove Theorem 2. Input is now an arbitrary SLP  $\mathcal{S} = \{\Sigma, \mathcal{V}, \mathcal{D}\}$  of size  $n$  generating  $T$ . Basically what we consider is to simulate  $\text{TtoG}$  on  $\mathcal{S}$ , namely, compute  $\text{TtoG}(T)$  without decompressing  $\mathcal{S}$  explicitly. The recompression technique is celebrated for doing this kind of tasks (actually this is where “recompression” is named after). In Section 4.1, we present an algorithm  $\text{SimTtoG}$  that simulates  $\text{TtoG}$  in  $O(n \lg^2(N/n))$  time and  $O(n + z \lg(N/z))$  space. In Section 4.2, we present how to modify  $\text{SimTtoG}$  to obtain an  $O(n \lg(N/n))$ -time and  $O(n + z \lg(N/z))$ -space construction of our  $\text{LCE}$  data structure.

### 4.1 SimTtoG: Simulating TtoG on CFGs

We present an algorithm  $\text{SimTtoG}$  to simulate  $\text{TtoG}$  on  $\mathcal{S}$ . To begin with, we compute the CFG  $\mathcal{S}_0 = \{\Sigma_0, \mathcal{V}, \mathcal{D}_0\}$  obtained by replacing, for all variables  $X \in \mathcal{V}$  with  $\mathcal{D}(X) \in \Sigma$ , every occurrence of  $X$  in the righthand sides of  $\mathcal{D}$  with the letter generating  $\mathcal{D}(X)$ . Note that  $\Sigma_0$  is the set of terminals of  $\mathcal{S}_0$ , and  $\mathcal{S}_0$  generates  $T_0$ .  $\text{SimTtoG}$  transforms level by level  $\mathcal{S}_0$  into CFGs,  $\mathcal{S}_1 = \{\Sigma_1, \mathcal{V}, \mathcal{D}_1\}, \mathcal{S}_2 = \{\Sigma_2, \mathcal{V}, \mathcal{D}_2\}, \dots, \mathcal{S}_h = \{\Sigma_h, \mathcal{V}, \mathcal{D}_h\}$ , where each  $\mathcal{S}_h$  generates  $T_h$ . Namely, compression from  $T_h$  to  $T_{h+1}$  is simulated on  $\mathcal{S}_h$ . We can correctly compute the letters in  $\hat{\Sigma}_{h+1}$  while modifying  $\mathcal{S}_h$  into  $\mathcal{S}_{h+1}$ , and hence, we get all the letters of  $\text{TtoG}(T)$  in the end. We note that new variables are never introduced and the modification is done by rewriting righthand sides of the original variables.

Here we introduce the special formation of the CFGs  $\mathcal{S}_h$  (it is a generalization of SLPs in a different sense from RLSLPs): For any  $X \in \mathcal{V}$ ,  $\mathcal{D}_h(X)$  consists of an “arbitrary number” of letters and at most “two” variables. More precisely, the following condition holds:

For any variable  $X \in \mathcal{V}$  with  $\mathcal{D}(X) = \dot{X}\dot{X}$ ,  $\mathcal{D}_h(X)$  is either  $w_1\dot{X}w_2\dot{X}w_3$ ,  $w_1\dot{X}w_2$ ,  $w_2\dot{X}w_3$  or  $w_2$  with  $w_1, w_2, w_3 \in \Sigma_h^*$ , where  $w_1 = w_3 = \varepsilon$  if  $X$  is not the starting variable.

As opposed to SLPs and RLSLPs, we define the size of  $\mathcal{S}_h$  by the total lengths of righthand sides and denote it by  $|\mathcal{S}_h|$ .

#### 4.1.1 PComp on CFGs

We firstly show that the adjacency list of  $T_h$  can be computed efficiently.

**Lemma 16** (Lemma 6.1 of [9]). *For any odd  $h$  ( $0 \leq h < \hat{h}$ ), the adjacency list of  $T_h$ , whose size is  $O(|\mathcal{S}_h|)$ , can be computed in  $O(|\mathcal{S}_h| + n)$  time and space.*

*Proof.* For any variable  $X \in \mathcal{V}$ , let  $\text{VOcc}(X)$  denote the number of occurrences of the nodes labeled with  $X$  in the derivation tree of  $\mathcal{S}$ . It is well known that  $\text{VOcc}(X)$  for all variables can be computed in  $O(n)$  time and space on the DAG representation of the tree.<sup>5</sup> Also, for any variable  $X \in \mathcal{V}$ , let  $\text{LML}(X)$  and  $\text{RML}(X)$  denote the leftmost letter and respectively rightmost letter of  $\text{val}_{\mathcal{S}_h}(X)$ . We can compute  $\text{LML}(X)$  for all variables in  $O(|\mathcal{S}_h|)$  time by a bottom up computation, i.e.,  $\text{LML}(X) = \text{LML}(Y)$  if  $\mathcal{D}_h(X)$  starts with a variable  $Y$ , and  $\text{LML}(X) = w[1]$  if  $\mathcal{D}_h(X)$  starts with a non-empty string  $w$ . In a completely symmetric way  $\text{RML}(X)$  can be computed in  $O(|\mathcal{S}_h|)$  time.

Now observe that any occurrence  $i$  of a pair  $\hat{c}\check{c}$  in  $T_h$  can be uniquely associated with a variable  $X$  that is the label of the lowest node covering the interval  $[i..i+1]$  in the derivation tree of  $\mathcal{S}_h$  (recall that  $\mathcal{S}_h$  generates  $T_h$ ). We intend to count all the occurrences of pairs associated with  $X$  in  $\mathcal{D}_h(X)$ . For example, let  $\mathcal{D}_h(X) = \hat{X}w_2\check{X}$  with  $w_2 \in \Sigma_h^*$ . Then  $\hat{c}\check{c}$  appears *explicitly* in  $w_2$  or *crosses* the boundaries of  $\hat{X}$  and/or  $\check{X}$ . If  $\hat{c}\check{c}$  crosses the boundary of  $\hat{X}$ ,  $\text{RML}(\hat{X})$  is  $\hat{c}$  and  $\check{c}$  follows, i.e.,  $w_2[1] = \check{c}$  or  $w_2 = \varepsilon$  and  $\text{LML}(\check{X}) = \check{c}$ . Using  $\text{RML}(\hat{X})$  and  $\text{LML}(\check{X})$ , we can list in  $O(|\mathcal{D}_h(X)|)$  time and space all the explicit and crossing pairs in  $\mathcal{D}_h(X)$ . Using  $\text{RML}(\hat{X})$  and  $\text{LML}(\check{X})$ , we can compute in  $O(|\mathcal{D}_h(X)|)$  time and space a  $|\mathcal{D}_h(X)| - 1$  size multiset that lists all the explicit and crossing pairs in  $\mathcal{D}_h(X)$ . Each pair  $\hat{c}\check{c}$  with  $\hat{c} > \check{c}$  (resp.  $\hat{c} < \check{c}$ ) is listed by a quadruple  $(\hat{c}, \check{c}, 0, \text{VOcc}(X))$  (resp.  $(\check{c}, \hat{c}, 1, \text{VOcc}(X))$ ).  $\text{VOcc}(X)$  means that the pair has a weight  $\text{VOcc}(X)$  because the pair appears every time a node labeled with  $X$  appears in the derivation tree.

We compute such a multiset for every variable, which takes  $O(|\mathcal{S}_h|)$  time and space in total. Next we sort the obtained list in increasing order of the first three integers in a quadruple. Note that the maximum value of letters is  $O(z \lg(N/z))$  due to Lemma 14, and  $O(z \lg(N/z)) = O(n^2)$  since  $z \leq n$  and  $\lg N \leq n$  hold. Thus the sorting can be done in  $O(n)$  time and space by radix sort. Finally we can get the adjacency list of  $T_h$  by summing up weights of the same pair. The size of the list is clearly  $O(|\mathcal{S}_h|)$ .  $\square$

The next lemma shows how to implement PComp on CFGs:

**Lemma 17.** *For any odd  $h$  ( $0 \leq h < \hat{h}$ ), we can compute  $\mathcal{S}_{h+1}$  from  $\mathcal{S}_h$  in  $O(|\mathcal{S}_h| + n)$  time and space. In addition,  $|\mathcal{S}_{h+1}| \leq |\mathcal{S}_h| + 2n$ .*

*Proof.* We first compute the partition  $(\hat{\Sigma}_h, \check{\Sigma}_h)$  of  $\Sigma_h$ , which can be done in  $O(|\mathcal{S}_h| + n)$  time and space by Lemmas 16 and 8.

Given  $(\hat{\Sigma}_h, \check{\Sigma}_h)$ , we can detect all the positions of the pairs from  $\hat{\Sigma}_h\check{\Sigma}_h$  in the righthands of  $\mathcal{D}_h$ , which should be compressed. Some of the appearances of the pairs are explicit and the others are crossing. While explicit pairs can be compressed easily, crossing pairs need an additional treatment. In order to deal with crossing pairs, we first *uncross* them by popping out  $\text{LML}(Y)$  (resp.  $\text{RML}(Y)$ ) from  $\text{val}_{\mathcal{S}_h}(Y)$  iff  $\text{LML}(Y) \in \check{\Sigma}_h$  (resp.  $\text{RML}(Y) \in \hat{\Sigma}_h$ ) for every variable  $Y$  other than the starting variable. More precisely, we do the following:

**PopInLet** For any variable  $X$ , if  $\mathcal{D}_h(X)[i] = Y \in \mathcal{V}$  with  $i > 1$  ( $i \geq 1$  if  $X$  is the starting variable) and  $\text{LML}(Y) \in \check{\Sigma}_h$ , replace the occurrence of  $Y$  with  $\text{LML}(Y)Y$ ; if  $\mathcal{D}_h(X)[i] = Y \in \mathcal{V}$  with  $i < |\mathcal{D}_h(X)|$  ( $i \leq |\mathcal{D}_h(X)|$  if  $X$  is the starting variable) and  $\text{RML}(Y) \in \hat{\Sigma}_h$ , replace the occurrence of  $Y$  with  $Y\text{RML}(Y)$ .

**PopOutLet** For any variable  $X$  other than the starting variable, if  $\mathcal{D}_h(X)[1] \in \check{\Sigma}_h$ , remove the first letter of  $\mathcal{D}_h(X)$ ; and if  $\mathcal{D}_h(X)[|\mathcal{D}_h(X)|] \in \hat{\Sigma}_h$ , remove the last letter of  $\mathcal{D}_h(X)$ .<sup>6</sup>

**PopOutLet** removes  $\text{LML}(Y)$  and  $\text{RML}(Y)$  from  $\text{val}_{\mathcal{S}_h}(Y)$  if they can be a part of a crossing pair and **PopInLet** introduces the removed letters into appropriate positions in  $\mathcal{D}_h$  so that the modified  $\mathcal{S}_h$  keeps to generate  $T_h$ . Notice that for each variable  $X$  the positions where letters popped-in is at most two (four if

<sup>5</sup>It is enough to compute  $\text{VOcc}(X)$  once at the very beginning of **SimTtoG**.

<sup>6</sup>If  $X$  becomes empty, we will clear all the appearance of  $X$  in  $\mathcal{D}_h$ .



$X$  is the starting variable) and there is at least one variable that has no variables below, and hence, the size of  $\mathcal{S}_h$  increases at most  $2n$ . The uncrossing can be conducted in  $O(|\mathcal{S}_h| + n)$  time.

Since all the pairs to be compressed become explicit now, we can easily conduct **BComp** in  $O(|\mathcal{S}_h| + n)$  time. We scan righthand sides in  $O(|\mathcal{S}_h|)$  time and list all the occurrences of pairs to be compressed. Each occurrence of pair  $\acute{c}\grave{c} \in \acute{\Sigma}\grave{\Sigma}$  is listed by a triple  $(\acute{c}, \grave{c}, p)$ , where  $p$  is the pointer to the occurrence. Then we sort the list according to the pair of integers  $(\acute{c}, \grave{c})$ , which can be done in  $O(|\mathcal{S}_h| + n)$  time and space by radix sort because  $\acute{c}$  and  $\grave{c}$  are  $O(n^2)$ . Finally, we replace each pair at position  $p$  with a fresh letter based on the rank of  $(\acute{c}, \grave{c})$ .  $\square$

#### 4.1.2 BComp on CFGs

For any even  $h$  ( $0 \leq h < \hat{h}$ ), **BComp** can be implemented in a similar way to **PComp** of Lemma 17. A block  $T_h[b..e]$  of length  $\geq 2$  is uniquely associated with a variable  $X$  that is the label of the lowest node covering the interval  $[b - 1..e + 1]$  in the derivation tree of  $\mathcal{S}_h$  (if  $b = 0$  or  $e = |T_h|$ , the block is associated with the starting variable). Note that we take  $[b - 1..e + 1]$  rather than  $[b..e]$  to be sure that the block cannot extend outside the variable. Some blocks are explicitly written in  $\mathcal{D}_h(X)$  and some others are crossing the boundaries of variables in  $\mathcal{D}_h(X)$ . The numbers of explicit blocks and crossing blocks in  $\mathcal{D}_h$  is at most  $|\mathcal{S}_h|$  and  $2n$ , respectively. The crossing blocks can be uncrossed in a similar way to uncrossing pairs. Then **BComp** can be done by replacing all the blocks with fresh letters on righthand sides of  $\mathcal{D}_h$ .

However here we have a problem. Recall that in order to give a unique letter to a block  $c^d$ , we have to sort the pairs of integers  $(c, d)$  (see Lemma 7). Since  $d$  might be exponentially larger than  $|\mathcal{S}_h| + n$ , radix sort cannot be executed in  $O(|\mathcal{S}_h| + n)$  time and space. In Section 6.2 of [9], Jeż showed how to solve this problem by tweaking the representation of lengths of long blocks, but its implementation and analysis are involved.<sup>7</sup>

We show in Lemma 18 our new observation, which leads to a simpler implementation and analysis of **BComp**. We say that a block  $c^d$  is *short* if  $d = O(|\mathcal{S}_h| + n)$  and *long* otherwise. Also, we say that a variable is *unary* iff its righthand side consists of a single block.

**Lemma 18.** *For any even  $h$  ( $0 \leq h < \hat{h}$ ), a block  $T_h[b..e] = c^d$  is short if it does not include a substring generated from a unary variable.*

*Proof.* Consider the derivation tree of  $\mathcal{S}_h$  and the shortest path from  $T_h[b]$  to  $T_h[e]$ . Let  $X_1 X_2 \cdots X_{m'} \cdots X_m$  be the sequence of labels of internal nodes on the path, where  $X_{m'}$  corresponds to the lowest common ancestor of  $T_h[b]$  and  $T_h[e]$ . Since SLPs have no loops in the derivation tree,  $X_1, \dots, X_{m'}$  are all distinct. Similarly  $X_{m'+1}, \dots, X_m$  are all distinct. Since a unary variable is not involved to generate the block, it is easy to see that  $d \leq \sum_{i=1}^m |\mathcal{D}_h(X_i)| \leq 2|\mathcal{S}_h|$  holds.  $\square$

Lemma 18 implies that most of blocks we find during the compression are short, which can be sorted efficiently by radix sort. If there is a long block in  $\mathcal{D}_h$ , an occurrence of a unary variable  $X$  must be involved to generate the block. Since **BComp** at level  $h$  pops out all the letters from  $X$  and removes the occurrences of  $X$  in  $\mathcal{D}_h$ , there are at most  $2n$  long blocks in total. The number of long blocks can also be upper bounded by  $2N/n$  with a different analysis based on the following fact:

**Fact 19.** *If a substring of original text  $T$  generated from a long block overlaps with that generated from another long block, one substring must include the other, and moreover, the shorter block is completely included in “one” letter of the longer block. Hence the length of the substring of the longer block is at least  $n$  times longer than that of the shorter block.*

Let us consider the long blocks that generate substrings whose lengths are  $[n^i..n^{i+1}]$  for a fixed integer  $i \geq 1$ . By Fact 19, the substrings cannot overlap, and hence, the number of such long blocks is at most  $N/n^i$ . Therefore, the total number of long blocks is at most  $\sum_{i \geq 1} N/n^i \leq 2N/n$ . Thus we get the following lemma.

**Lemma 20.** *There are at most  $O(\min(n, N/n))$  long blocks found during **SimTtoG**.*

By Lemma 20, we can employ a standard comparison-base sorting algorithm to sort all long blocks in  $O(n \lg(\min(n, N/n)))$  time in total. In particular, **BComp** of one level can be implemented in the following complexities:

<sup>7</sup>Note that Section 6.2 of [9] also takes care of the case where the word size is  $\Theta(\lg n)$  rather than  $\Theta(\lg N)$ . We do not consider the  $\Theta(\lg n)$ -bits model in this paper because using  $\Theta(\lg N)$  bits to store the length of string generated by every letter is crucial for extract and LCE queries. However, we believe that our new observation stated in Lemma 18 will simplify the analysis for the  $\Theta(\lg n)$ -bits model, too.

**Lemma 21.** For any even  $h$  ( $0 \leq h < \hat{h}$ ), we can compute  $\mathcal{S}_{h+1}$  from  $\mathcal{S}_h$  in  $O(|\mathcal{S}_h| + n + m \lg m)$  time and  $O(|\mathcal{S}_h| + n)$  space, where  $m$  is the number of long blocks in  $\mathcal{D}_h$ . In addition,  $|\mathcal{S}_{h+1}| \leq |\mathcal{S}_h| + 2n$ .

### 4.1.3 The complexities of SimTtoG

**Theorem 22.** SimTtoG runs in  $O(n \lg^2(N/n))$  time and  $O(n \lg(N/n))$  space.

*Proof.* Using PComp and BComp implemented on CFGs (see Lemma 17 and 21), SimTtoG transforms level by level  $\mathcal{S}_0$  into  $\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_{\hat{h}}$ . In each level, the size of CFGs can increase at most  $2n$  by the procedure of uncrossing. Since  $|\mathcal{S}_h| = O(n \lg N)$  for any  $h$  ( $0 \leq h < \hat{h}$ ), we get the time complexity  $O(n \lg^2 N)$  by simply applying Lemmas 17 and 21.

We can improve it to  $O(n \lg^2(N/n))$  by a similar trick used in the proof of Lemma 14. At some level  $h'$  where  $|T_{h'}|$  becomes less than  $n$ , we decompress  $\mathcal{S}_{h'}$  and switch to TtoG, which transforms  $T_{h'}$  into  $T_{\hat{h}}$  in  $O(n)$  time by Lemma 10. We apply Lemmas 17 and 21 only for  $h$  with  $0 \leq h < h'$ . Since  $h' = O(\lg(N/n))$ ,  $|\mathcal{S}_h| = O(n \lg(N/n))$  for any  $h$  ( $0 \leq h < h'$ ). Therefore, we get the time complexity  $O(n \lg^2(N/n))$ . The space complexity is bounded by the maximum size of CFGs  $\mathcal{S}_0, \mathcal{S}_1, \dots, \mathcal{S}_{h'}$ , which is  $O(n \lg(N/n))$ .  $\square$

## 4.2 GtoG: $O(n \lg(N/n))$ -time recompression

We modify SimTtoG slightly to run in  $O(n \lg(N/n))$  time and  $O(n + z \lg(N/z))$  space. The idea is the same as what has been presented in Section 6.1 of [9]. The problem of SimTtoG is that the sizes of intermediate CFGs  $\mathcal{S}_h$  can grow up to  $O(n \lg(N/n))$ . If we can keep their sizes to  $O(n)$ , everything goes fine. This can be achieved by using two different types of partitions of  $\Sigma_h$  for PComp: One is for compressing  $T_h$  by a constant factor, and the other for compressing  $|\mathcal{S}_h|$  by a constant factor (unless  $|\mathcal{S}_h|$  is too small to compress). Recall that the former partition has been used in TtoG and SimTtoG, and the partition is computed from the adjacency list of  $T_h$  by Algorithm 1. Algorithm 1 can be extended to work on a set of strings by just inputting the adjacency list from a set of strings. Then, we can compute the partition for compressing  $|\mathcal{S}_h|$  by a constant factor by considering the adjacency list from a set of strings in the righthand sides of  $\mathcal{D}_h$ . The adjacency list can be easily computed in  $O(|\mathcal{S}_h| + n)$  time and space by modifying the algorithm described in the proof of Lemma 16: We just ignore the weight  $\text{VOcc}(X)$ , i.e., use a unit weight 1 for every listed pair. Using the two types of partitions alternately, we can compress strings by a constant factor while keeping the intermediate CFGs to  $O(n)$ .

We denote the modified algorithm by GtoG and the resulting RLSLP by  $\text{GtoG}(\mathcal{S})$ . Note that  $\text{GtoG}(\mathcal{S})$  is not identical to TtoG( $T$ ) in general because the partitions used in GtoG change depending on the input  $\mathcal{S}$ . Still the height of  $\text{GtoG}(\mathcal{S})$  is  $O(\lg N)$  and the properties of *PSeqs* hold. Hence we can support LCE queries on  $\text{GtoG}(\mathcal{S})$  as we did on TtoG( $T$ ) by Lemma 15.

## 4.3 Proof of Theorem 2

*Proof of Theorem 2.* Let  $\mathcal{S}$  be an input SLP of size  $n$  generating  $T$ . We compute  $\text{GtoG}(\mathcal{S})$  in  $O(n \lg(N/n))$  time and  $O(n + z \lg(N/z))$  space as described in Section 4.2. Since the height of TtoG( $T$ ) is  $O(\lg N)$ , we can support Extract( $i, \ell$ ) queries in  $O(\lg N + \ell)$  time due to Lemma 6.  $\text{GtoG}(\mathcal{S})$  supports LCE queries in  $O(\lg N)$  time in the same way as what was described in Lemma 15.  $\square$

**Acknowledgements.** The author was supported by JSPS KAKENHI Grant Number 16K16009.

## References

- [1] Michael A. Bender, Martin Farach-Colton, Giridhar Pemmasani, Steven Skiena, and Pavel Sumazin. Lowest common ancestors in trees and directed acyclic graphs. *J. Algorithms*, 57(2):75–94, 2005.
- [2] P. Bille, P. H. Cording, I. L. Gørtz, B. Sach, H. W. Vildhøj, and Søren Vind. Fingerprints in compressed strings. In *Proc. WADS 2013*, pages 146–157, 2013.
- [3] Philip Bille, Anders Roy Christiansen, Patrick Hagge Cording, and Inge Li Gørtz. Finger search in grammar-compressed strings. *arXiv:1507.02853v3*.
- [4] Philip Bille, Inge Li Gørtz, Mathias Bæk Tejs Knudsen, Moshe Lewenstein, and Hjalte Wedel Vildhøj. Longest common extensions in sublinear space. In *Proc. CPM 2015*, pages 65–76, 2015.
- [5] Philip Bille, Inge Li Gørtz, Benjamin Sach, and Hjalte Wedel Vildhøj. Time-space trade-offs for longest common extensions. *J. Discrete Algorithms*, 25:42–50, 2014.
- [6] Dan Gusfield. *Algorithms on Strings, Trees, and Sequences*. Cambridge University Press, 1997.
- [7] Tomohiro I, Wataru Matsubara, Kouji Shimohira, Shunsuke Inenaga, Hideo Bannai, Masayuki Takeda, Kazuyuki Narisawa, and Ayumi Shinohara. Detecting regularities on grammar-compressed strings. *Inf. Comput.*, 240:74–89, 2015.
- [8] Artur Jeż. Approximation of grammar-based compression via recompression. *Theor. Comput. Sci.*, 592:115–134, 2015.
- [9] Artur Jeż. Faster fully compressed pattern matching by recompression. *ACM Transactions on Algorithms*, 11(3):20:1–20:43, 2015.
- [10] Artur Jeż. One-variable word equations in linear time. *Algorithmica*, 74(1):1–48, 2016.
- [11] Artur Jeż. Recompression: A simple and powerful technique for word equations. *J. ACM*, 63(1):4, 2016.
- [12] Shirou Maruyama, Masaya Nakahara, Naoya Kishiue, and Hiroshi Sakamoto. Esp-index: A compressed index based on edit-sensitive parsing. *J. Discrete Algorithms*, 18:100–112, 2013.
- [13] Kurt Mehlhorn, R. Sundar, and Christian Urig. Maintaining dynamic sequences under equality tests in polylogarithmic time. *Algorithmica*, 17(2):183–198, 1997.
- [14] Takaaki Nishimoto, Tomohiro I, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda. Dynamic index and LZ factorization in compressed space. In *Proceedings of the Prague Stringology Conference 2016, Prague, Czech Republic, August 29-31, 2015*, pages 158–170, 2016.
- [15] Takaaki Nishimoto, Tomohiro I, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda. Fully dynamic data structure for LCE queries in compressed space. In *41st International Symposium on Mathematical Foundations of Computer Science, MFCS 2016, August 22-26, 2016 - Kraków, Poland*, pages 72:1–72:15, 2016.
- [16] Nicola Prezza. In-place longest common extensions. *arXiv:1608.05100v6*, 2016.
- [17] Simon J. Puglisi and Andrew Turpin. Space-time tradeoffs for longest-common-prefix array computation. In *Proc. ISAAC '08*, volume 5369 of *Lecture Notes in Computer Science*, pages 124–135. Springer, 2008.
- [18] Wojciech Rytter. Application of Lempel-Ziv factorization to the approximation of grammar-based compression. *Theoretical Computer Science*, 302(1–3):211–222, 2003.
- [19] Yuka Tanimura, Tomohiro I, Hideo Bannai, Shunsuke Inenaga, Simon J. Puglisi, and Masayuki Takeda. Deterministic sub-linear space LCE data structures with efficient construction. In *27th Annual Symposium on Combinatorial Pattern Matching, CPM 2016, June 27-29, 2016, Tel Aviv, Israel*, pages 1:1–1:10, 2016.