

Generalized Substring Compression

Orgad Keller^a, Tsvi Kopelowitz^a, Shir Landau Feibish^{b,1}, Moshe Lewenstein^{a,2}

^a*Department of Computer Science, Bar-Ilan University, Ramat-Gan, Israel*

^b*Blavatnik School of Computer Sciences, Tel-Aviv University, Israel*

Abstract

In *substring compression* one is given a text to preprocess so that, upon request, a compressed substring is returned. *Generalized substring compression* is the same with the following twist. The queries contain an additional context substring (or a collection of context substrings) and the answers are the substring in compressed format, where the context substring is used to make the compression more efficient.

We focus our attention on *generalized substring compression* and present the first non-trivial *correct* algorithm for this problem. Inherent to our algorithm is a new method for finding the *bounded longest common prefix* of substrings, which may be of independent interest. In addition, we propose an efficient algorithm for substring compression which makes use of *range successor* queries.

We present several tradeoffs for both problems. For compressing the substring $S[i..j]$ (possibly with the substring $S[\alpha.. \beta]$ as a context), the best query times we achieve are $O(C)$ and $O(C \log(\frac{j-i}{C}))$ for substring compression query and generalized substring compression query, respectively, where C is the number of phrases encoded.

A preliminary version of this paper has been presented in [21].

Keywords: Data compression, Lempel-Ziv compression, Suffix tree, Range searching

1. Introduction

While string compression has been studied for decades, substring compression is relatively lightly studied. The topic was introduced in [8], where a set of problems concerning substring compression focusing on the compression algorithm of Lempel and Ziv [35] was presented. They deal mainly with two

Email addresses: kellero@cs.biu.ac.il (Orgad Keller), kopelot@cs.biu.ac.il (Tsvi Kopelowitz), lfshir@gmail.com (Shir Landau Feibish), moshe@cs.biu.ac.il (Moshe Lewenstein)

¹The author participated in this research when she was a student at Bar-Ilan University.

²The research of the author was supported by a BSF grant 2010437, a Google Research Award and by a GIF grant 1147/2011.

variants of this topic, namely, given a string, what is the compressibility of different substrings of that string, both in the sense of the actual compression of the substrings and in the sense of comparing which of the substrings is the least or most compressible.

The goal of our research is to find the inherent connection between the compressed representation of a string and that of its substrings. Such a connection could have interesting practical applications, and may lead the way to finding further connections between certain properties of strings and that of their substrings. In addition, it may be interesting to further investigate the topic of substring compression using other compression methods.

We address the following problems: in the *substring compression query* (SCQ) problem, we wish to compress a given substring of the string S , denoted by start and end location. Note that we are allowed to preprocess S beforehand, so that we would be able to answer this query for any substring in S without having to scan the substring during query time. In its generalized and more powerful version, the *generalized substring compression query* (GSCQ) problem, we wish to compress the substring according to a given context taken from S as well. In both problems, our goal is to provide query times which are proportional to the size of the *compressed* substring as opposed to the size of the substring in its non-compressed form.

1.0.1. Applications

The issue of substring compression has interesting implications for a variety of practical applications. Recent works such as those presented in [7, 9, 31] for example, use compression of biological sequences as a basis of comparison between different sequences, and their information content. Compression of subsequences can therefore be used to perform such comparisons in a more efficient and accurate manner.

The result presented in [18], uses straight-line-program compression in order to speed up computation of edit-distance. This result relies heavily on the findings of Rytter [32].³ There, Rytter proved that an LZ77 [35] encoding can be transformed to a straight-line program quickly and without large expansion. Therefore, one may be able to use substring compression to speed up the edit distance computation of substrings, which may be a problem of independent interest.

1.1. Our Results

1. Our main result is an efficient and innovative algorithm for the *generalized substring compression query*, introduced in [8]. There an algorithm was suggested, however it is incorrect [29]: it overlooked the inherent added difficulty of the generalized problem which uses a bounded context, dismissing it as trivial, while it is in fact the essence of the generalized problem. The additional bounded context requires a different algorithm

³See also [6] where the same problem was independently addressed by Charikar et al.

Table 1: Results

Prob.	Query Time	Space	Source
GSCQ	$O(C_{\alpha,\beta}(i, j) \log (\frac{j-i}{C_{\alpha,\beta}(i,j)}))$	$O(n^{1+\epsilon})$	new
	$O(C_{\alpha,\beta}(i, j) \log (\frac{j-i}{C_{\alpha,\beta}(i,j)}) \log \log n)$	$O(n \log^\epsilon n)$	new
	$O(C_{\alpha,\beta}(i, j) (\log (\frac{j-i}{C_{\alpha,\beta}(i,j)}) \log \log n + (\log \log n)^2))$	$O(n \log \log n)$	new
	$O(C_{\alpha,\beta}(i, j) \log (\frac{j-i}{C_{\alpha,\beta}(i,j)}) \log^\epsilon n)$	$O(n)$	new
SCQ	$O(C(i, j))$	$O(n^{1+\epsilon})$	new
	$O(C(i, j) \log \log n)$	$O(n \log^\epsilon n)$	new
	$O(C(i, j) (\log \log n)^2)$	$O(n \log \log n)$	new
	$O(C(i, j) \log^\epsilon n)$	$O(n)$	new
	$O(C(i, j) \log n \log \log n)$	$O(n \log^\epsilon n)$	[8]

than the context within the substring to be compressed, that we will describe in detail in this paper. Therefore, the solution provided in [8] in fact does not solve the problem. Our solution for this problem is based on a solution to finding the *bounded longest common prefix* (BLCP) of two substrings, which is a notion we will introduce shortly.

2. In addition, we improve results shown for the *substring compression query*. Our result is based mainly on an improved solution for finding the *interval longest common prefix* (ILCP) of two substrings. This is done using an efficient solution for the problem of *range successor* [25]⁴, and not on the more classical *range reporting* problem (see, for instance [1, 5]), used by [8] and numerous other indexing-related papers [14, 2, 15, 27]. This constitutes a different method in order to reduce the substring compression query problem to the geometric problem. See [26] for a survey on the connection between text indexing and various range searching techniques.

Our solutions are based on a variety of tools, such as suffix trees, lowest common ancestor queries, level ancestor queries, and several kinds of range searching structures. As a result, solutions to both SCQ and GSCQ constitute tradeoffs between query times and space, due to the choice of range searching structures to be used. Denote C as the number of phrases encoded. A comparison of the results is presented in Table 1.

Note that range problems on strings has garnered much interest lately with different papers exploring different aspects of range problems. We refer the interested reader to [8, 21, 22, 30, 27, 4, 19, 23].

The rest of our paper is organized as follows: in Section 2, we give some preliminaries and problem definitions. In Sections 3.3 and 3.2, we describe our solutions for finding the BLCP and ILCP. In Section 4, we present the outline of the query algorithm’s main loop, which is roughly common to both the SCQ and GSCQ problems. In Sections 5 and 6, we present the solutions and analysis

⁴The range successor problem was introduced in [25] under the name *range searching for minimum*. The name “range successor” is used in [22, 30]. An almost-identical problem is the *range next value* problem [12] that will be discussed later.

for SCQ and GSCQ.

2. Problem Definitions and Preliminaries

2.1. Preliminary Definitions and Notations

Given a string S , $|S|$ is the length of S . Throughout this paper we denote $n = |S|$. An integer i is a *location* or a *position* in S if $i = 1, \dots, |S|$. The substring $S[i..j]$ of S , for any two positions $i \leq j$, is the substring of S that begins at index i and ends at index j . Concatenation is denoted by juxtaposition. The *suffix* S_i of S is the substring $S[i..n]$.

The *suffix tree* [34, 33, 13, 28] of a string S , denoted $ST(S)$, is a compact trie of all the suffixes of $S\$$ (i.e., S concatenated with a delimiter symbol $\$ \notin \Sigma$, where Σ is the alphabet set). Each of its edges is labeled with a substring of S (actually, a representation of it, e.g., the start location and its length). The “compact” property is achieved by contracting nodes having a single child. The children of every node are sorted in the lexicographical order of the substrings on the edges leading to them. Consequently, each leaf of the suffix tree represents a suffix of S , and the leaves are sorted from left to right in the lexicographical order of the suffixes that they represent. $ST(S)$ requires $O(n)$ space. Algorithms for the construction of a suffix tree enable $O(n)$ preprocessing time when $|\Sigma|$ is constant, and $O(n \log \min(n, |\Sigma|))$ time when $|\Sigma|$ is not. In fact, the suffix tree can be constructed in linear time even for alphabets drawn from a polynomially-sized range, see [13].

In addition, our algorithms make use of elements from the field of computational geometry; let $P = \{(x_1, y_1), \dots, (x_n, y_n)\}$ be a set of n points on an $[n] \times [n]$ grid. The following range searching query types are defined on P , for various types of a two-dimensional range R :

$\text{rangesucc}_y(R = [x, x'] \times [y, \infty])$: reports the single point of P that is included in the range and has a minimal y -coordinate, i.e., the point $\arg \min_{(x,y) \in P \cap R} y$. In other words, the resulting point is the successor to y (on the y -axis) having x -coordinate in the range $[x, x']$. The query types rangesucc_x and rangepred_x for successor and predecessor respectively, this time on the x -axis, are defined in the same manner.

$\text{emptiness}(R = [x, x'] \times [y, y'])$: returns “true” iff $P \cap R = \emptyset$.

2.1.1. An Overview of the Lempel-Ziv Algorithm

The LZ77 variation of the Lempel-Ziv algorithm works as follows: given an input string S of length n , the algorithm encodes the string in a greedy manner from left to right. At each step of the algorithm, suppose we have already encoded $S[1..k-1]$, we search for the location t , such that $1 \leq t \leq k-1$, for which the longest common prefix of $S[k..n]$ and the suffix S_t , is maximal. Once we have found the desired location, suppose the aforementioned longest common prefix is the substring $S[t..r]$, a phrase will be added to the output which will include the encoding of the distance to the substring (i.e., the value

$k - t$) and the length of the substring, (i.e., the value $r - t + 1$). The algorithm continues by encoding $S[k + (r - t + 1) .. n]$. Finally, we denote the output of the LZ77 algorithm on the input S as $\text{LZ}(S)$.

The string S may be encoded within the context of the string T . We denote this by $\text{LZ}(S \mid T)$. The practical meaning of this is that the result is as if the algorithm was performed on the concatenated string $T\$S$, where $\$$ is a symbol that does not appear in neither S nor T , however, only the portion of $\text{LZ}(T\$S)$ which represents the compression of S is output by the algorithm. Some exceptions apply to this rule as will be described later. An example for the use of an additional context for compression and retrieval of genomic sequences can be seen in [24].

Recently, several works [10, 20, 16] have been done concerning a related problem, called *Lempel-Ziv factorization*. The Lempel-Ziv factorization of a string S is defined to be the decomposition $S = s_1 s_2 \dots s_v$ such that for each $i = 1, \dots, v - 1$, s_i is the longest prefix of $s_i s_{i+1} \dots s_v$ that appears in $s_1 \dots s_{i-1}$. If this prefix is empty, s_i will be a single character [11]. It may be interesting to see how this problem and the solutions proposed for it can be adapted to the problems we present in this work.

2.2. Problem Definitions

Given a string S of length n , we wish to preprocess S in such a way that allows us to efficiently answer the following queries:

Substring Compression Query (SCQ(i, j)): given any two indices i and j , such that $1 \leq i \leq j \leq n$, we wish to output $\text{LZ}(S[i .. j])$.

Generalized Substring Compression Query (GSCQ(i, j, α, β)): given any four indices i, j, α , and β , such that $1 \leq i \leq j \leq n$ and $1 \leq \alpha \leq \beta \leq n$, we wish to output $\text{LZ}(S[i .. j] \mid S[\alpha .. \beta])$.

Query times for both of the above query types will be strongly dependent on the number of phrases actually encoded. We denote these as $C(i, j)$ and $C_{\alpha, \beta}(i, j)$ for SCQ and GSCQ, respectively. Our results will rely on the two following primitives:

Bounded Longest Common Prefix (BLCP(k, l, r)): given k , and given positions l and r which induce the context substring $S[l .. r]$, we look for the longest common prefix of $S[k .. j]$ and a substring which starts at some location $l \leq t \leq r$ within the context. The substring chosen must not exceed the end of context. In other words, it must be a prefix of some substring $S[t .. r]$.

Interval Longest Common Prefix (ILCP(k, l, r)): given k, l, r , this time we look for the longest common prefix of $S[k .. j]$ and a substring which starts at some location $l \leq t \leq r$, without further constraints.

While it may not seem so at first glance, BLCP queries are more difficult to implement than ILCP. For example, consider two suffixes S_{t_1} and S_{t_2} , such that $l \leq t_1 < t_2 \leq r$, for which $|\text{LCP}(S_k, S_{t_1})| < |\text{LCP}(S_k, S_{t_2})|$ (where $\text{LCP}(S_1, S_2)$, for two strings S_1 and S_2 , stands for the *longest common prefix* of S_1 and S_2).

Some portion of the last characters of $\text{LCP}(S_k, S_{t_2})$ may not be eligible for consideration. Namely, if $|\text{LCP}(S_k, S_{t_2})|$ exceeds $r - t_2 + 1$ characters, $\text{LCP}(S_k, S_{t_2})$ exceeds location r , and therefore literally “grows out of context”. In that case, it may be that S_{t_1} will eventually be the suffix to be preferred. One should take into account that such a cut-off may pertain to $\text{LCP}(S_k, S_{t_1})$ as well. (Note: in the case $i - 1 \leq r < j$, if desired, one can allow a substring taken from the context to exceed r . This is a trivial extension to the algorithm for ILCP.)

3. Limiting the Longest Common Prefix: Answering ILCP and BLCP Queries

3.1. Preprocessing Motivation

We begin the preprocessing by constructing the suffix tree of S , $\text{ST}(S)$. In the suffix tree, each leaf ℓ is associated with a suffix of $S\$$ and is therefore marked with an integer $y(\ell)$ which is the start location of that suffix. We also mark each leaf ℓ with an integer $x(\ell)$ which is the lexicographical rank of the suffix associated with ℓ within the set of all suffixes of T (this is done by using one depth-first traversal, in which we number the leaves from left to right). We then preprocess the set $P = \{(x(\ell), y(\ell)) \mid \ell \text{ is a leaf in } \text{ST}(S)\} \subseteq [n + 1]^2$ for the range searching query types mentioned before. We will refer to the points in P as *suffix points*.

Suppose we search $\text{ST}(S)$ for some substring $S[l..r]$, we can find all the occurrence positions of $S[l..r]$ in S , by traversing $\text{ST}(S)$ from the root downwards according to the symbols in $S[l..r]$, until either (1) the next symbol of the pattern cannot be found at our current location in the tree — in this case we conclude that $S[l..r]$ does not occur in S ; (2) the pattern is exhausted and we conclude the traversal at a node v in $\text{ST}(S)$ (or the edge leading to it from its parent, for that matter). Let v be the node in which the search ended. All the leaves in the subtree rooted at v , denoted T_v , correspond to occurrences of $S[l..r]$ in S . Hence the set $Y_v = \{y(\ell) \mid \ell \text{ is a leaf in } T_v\}$ is the set of all occurrence positions of $S[l..r]$ in S . From the properties of the suffix tree it follows that the set $X_v = \{x(\ell) \mid \ell \text{ is a leaf in } T_v\}$ forms a consecutive range of values in $[n + 1]$. This is exactly the range $X_v = [x(l_v), x(r_v)]$, where l_v and r_v are the leftmost and rightmost leaves in T_v , respectively. It therefore holds that for a leaf ℓ , ℓ is a leaf in T_v iff $x(\ell) \in [x(l_v), x(r_v)]$. In other words: $x(\ell) \in [x(l_v), x(r_v)]$ iff $S[l..r]$ appears in S at location $y(\ell)$.

Notice that each node u in the suffix tree has two different notions of *depth*: the ordinary perception of depth of a node in a tree, denoted $\text{depth}(u)$, and the length of the string u represents (derived by the concatenation of edge labels on the path from the root to u), denoted $\text{length}(u)$. Now let S_i and S_j be two suffixes of S , and consider the *longest common prefix* of S_i and S_j , denoted $\text{LCP}(S_i, S_j)$. Let ℓ_i and ℓ_j be the leaves corresponding to S_i and S_j , respectively (i.e., $i = y(\ell_i)$ and $j = y(\ell_j)$). Then $|\text{LCP}(S_i, S_j)| = \text{length}(\text{LCA}(\ell_i, \ell_j))$, where $\text{LCA}(\ell_i, \ell_j)$ is the *lowest common ancestor* of ℓ_i and ℓ_j .

3.2. Answering ILCP queries

Here our primary goal is to obtain an efficient way of finding $\text{ILCP}(k, l, r)$, that is, given k, l, r , we seek the longest common prefix of $S[k..j]$ and a substring which starts at some location $l \leq t \leq r$. An example of this constraint can be seen in Fig. 1.

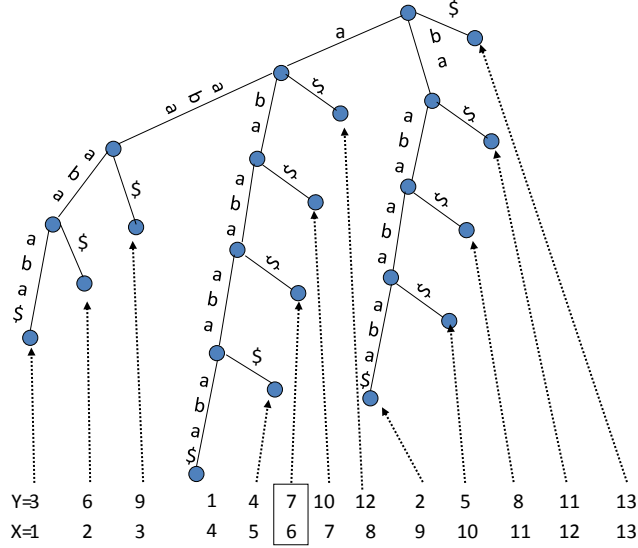


Figure 1: The suffix tree for the string $S = abaabaabaaba\$$ is shown. The x values indicate the lexicographical rank of the suffix. The y values indicate the index of the suffix in S . Assume we are compressing the substring: $S[4..9]$ that is $abaaba$, and that, so far, we have encoded the substring: $S[4..6]$, meaning $abaaba$. The marked suffix indicates the suffix that should be encoded next.

Recall that when searching for $\text{ILCP}(k, l, r)$, while the resulting substring must start at some location $l \leq t \leq r$, it is allowed to exceed location r . This is the equivalent to finding the location $l \leq t \leq r$, for which the longest common prefix of $S[k..j]$ and the suffix S_t , is maximal.

Consider the suffix S_k . Clearly, it is sufficient to find the location $t \in [l, r]$ for which $|\text{LCP}(S_k, S_t)|$ is maximized, i.e., $t = \arg \max_{z \in [l, r]} |\text{LCP}(S_k, S_z)|$ (without necessarily computing the value $|\text{LCP}(S_k, S_t)|$ at this stage). Once the aforementioned location t is found, we compute $|\text{LCP}(S_k, S_t)|$. Therefore, to summarize, we have two steps: (1) finding the location t , and (2) computing $|\text{LCP}(S_k, S_t)|$.

3.2.1. Finding the Start Location t

We use a reduction to the problem of range successor on a grid, as opposed to the *range reporting* used in [8]. An example of the geometric representation of the scenario depicted in Fig. 1 can be seen in Fig. 2.

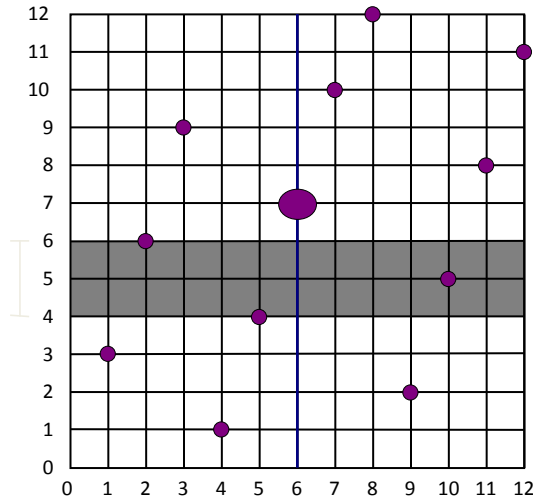


Figure 2: The grid depicts the geometric representation of the suffix tree for the string $S = abaabaabaaba\$$. The enlarged point at $(6,7)$ represents the suffix S_7 , as the substring S_7 , yet to be encoded, starts at position 7 in the string S . The grayed area of the grid represents the part of the substring which has already been encoded, that is $S[4..6]$. When finding the start location t , we will be limited to using points found in the gray area.

Consider the suffix S_k , and consider the set of suffixes $\Gamma = \{S_l, \dots, S_r\}$. Since $|\text{LCP}(S_k, S_t)| = \max_{z \in [l, r]} |\text{LCP}(S_k, S_z)|$, S_t is in fact the suffix lexicographically closest to S_k , out of all the suffixes of the set Γ .

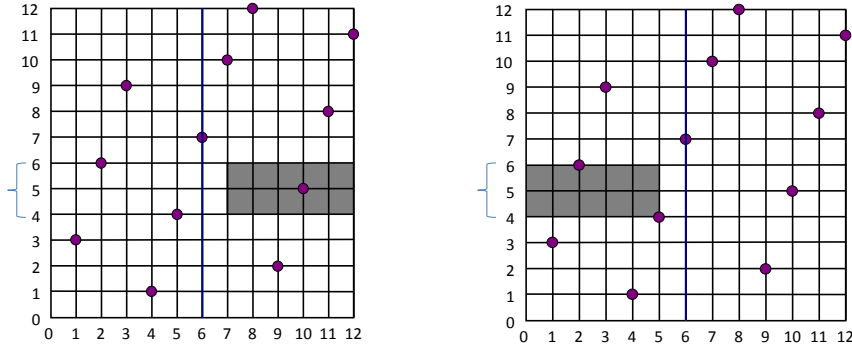
We will first assume that we are searching for a suffix S_{t_1} , such that the suffix S_{t_1} is *lexicographically smaller* than S_k . The process for the case where the suffix chosen is lexicographically greater than S_k is symmetric. Therefore, all we are required is to choose the best of both, i.e., the option yielding the greater $|\text{LCP}(S_k, S_t)|$ value.

Since we have assumed w.l.o.g. that S_{t_1} is lexicographically smaller than S_k , we have actually assumed that $x(\ell_{t_1}) < x(\ell_k)$, or equivalently, that ℓ_{t_1} appears to the *left* of ℓ_k in the suffix tree. Incorporating the lexicographical ranks of S_k and S_{t_1} into the expression, t_1 is actually the value which maximizes the expression $\max\{x(\ell_{t_1}) \mid l \leq t_1 \leq r \text{ and } x(\ell_{t_1}) < x(\ell_k)\}$. Notice that $t_1 = y(\ell_{t_1})$.

Now consider the set $P = \{(x(\ell), y(\ell)) \mid \ell \text{ is a leaf in } \text{ST}(S)\}$. Assuming indeed $x(\ell_{t_1}) < x(\ell_k)$, we are interested in finding the maximal value $x(\ell_{t_1})$, such that $x(\ell_{t_1}) < x(\ell_k)$, and $l \leq y(\ell_{t_1}) \leq r$. It immediately follows that the point $(x(\ell_{t_1}), y(\ell_{t_1})) \in P$ is the suffix point in the range $[-\infty, x(\ell_k) - 1] \times [l, r]$ having

maximal x -coordinate, and therefore can be obtained efficiently by querying $\text{rangepred}_x([-\infty, x(\ell_x) - 1] \times [l, r])$. Once we have found the point $(x(\ell_{t_1}), y(\ell_{t_1}))$, we can locate ℓ_{t_1} , as it is the $x(\ell_{t_1})$ -th leaf from the left. The leaf ℓ_{t_1} will be of importance in the next section.

Equivalently, there exists t_2 such that S_{t_2} is the suffix lexicographically larger than S_k and closest to it. In other words, we assume $x(\ell_{t_2}) > x(\ell_k)$, or equivalently, that ℓ_{t_2} appears to the *right* of ℓ_k in the suffix tree. t_2 can be found using a symmetric procedure: we query for the $\text{rangesucc}_x([x(\ell_x) + 1, \infty] \times [l, r])$. An example of the queries performed can be seen in Fig. 3.



(a) $\text{rangesucc}_x([x(\ell_x) + 1, \infty] \times [l, r])$,

(b) $\text{rangepred}_x([-\infty, x(\ell_x) - 1] \times [l, r])$,

Output: the point (10, 5).

Output: the point (5, 4).

Figure 3: Each grid depicts the geometric representation of the suffix tree for the string $S = \text{abaabaabaaba}\$$. The values given for the example queries are: $x = 6$, $l = 4$ and $r = 6$. The chosen suffix in this case would be S_4 since $|\text{LCP}(S_7, S_4)|$ is greater than $|\text{LCP}(S_7, S_5)|$.

Determining whether $t = t_1$ or $t = t_2$ is implemented by calculating both $|\text{LCP}S_k, S_{t_1}|$ and $|\text{LCP}S_k, S_{t_2}|$, and choosing the larger of the two. The exact method of calculating both values is described next.

3.2.2. Computing $|\text{LCP}(S_k, S_t)|$

Consider ℓ_k and ℓ_t as described above. Since $|\text{LCP}(S_i, S_j)| = \text{length}(\text{LCA}(\ell_i, \ell_j))$ for any i and j , it is sufficient to find the node $w = \text{LCA}(\ell_k, \ell_t)$ and then to compute $\text{length}(w)$. Using the methods of Harel and Tarjan [17], an LCA query can be answered in constant time. If the value $\text{length}(u)$ for each node u has been stored in u beforehand, we conclude the value $\text{length}(w)$ is obtainable in $O(1)$ time.

3.3. Answering BLCP queries

Assume that we wish to answer the query $\text{BLCP}(k, l, r)$. Consider the suffix S_k represented by the path from the root to ℓ_k . For any other suffix S_t and an integer $d \geq 0$, let $\text{LCP}_d(S_k, S_t)$ be the longest common prefix of S_k and S_t , truncated to at most d characters (i.e., if $\text{LCP}(S_k, S_t)$ exceeds d characters, we will leave only the first d and discard the others). By the definition of $\text{BLCP}(k, l, r)$, we are limited to finding substrings that do not exceed location r . Therefore, we actually wish to find $t' = \arg \max_{t \in [l, r]} |\text{LCP}_{r-t+1}(S_k, S_t)|$. Here notice that the term $r - t + 1$ is the maximal length of the portion of S_t we can use, according to the constraints. Also notice that by definition, $|\text{LCP}_{r-t+1}(S_k, S_t)| = \min\{|\text{LCP}(S_k, S_t)|, r - t + 1\}$. If several positions t that maximize the above expression exist, we define t' to be the leftmost such position, i.e., the smallest such value.

Definitions and Notations. For a node u , let $\text{path}(u)$ be the path in $\text{ST}(T)$ from the root to u . With a slight abuse of notation, we will also use $\text{path}(u)$ to denote the set of nodes participating in such a path. For a suffix S_t we also define $\text{path}(S_t) = \text{path}(\ell_t)$ where ℓ_t is the leaf representing S_t in the suffix tree. In the context of a $\text{BLCP}(k, l, r)$ query, a suffix S_t is said to be *relevant* if $l \leq t \leq r$.

A suffix S_t is said to be *eligible* at a node $u \in \text{path}(S_k)$ if the string represented by $\text{path}(u)$ is a prefix of $\text{LCP}_{r-t+1}(S_k, S_t)$. Consider the suffix S_k represented by the path from the root to ℓ_k . As suffixes S_t with greater $|\text{LCP}(S_k, S_t)|$ values branch out of this path at a later stage (i.e., leave this path at nodes of greater depth), we are interested in suffixes which share a large portion of this path. However, as $|\text{LCP}_{r-t+1}(S_k, S_t)|$ (and not $|\text{LCP}(S_k, S_t)|$) is the expression to be maximized, we are restricted by the eligibility of suffixes along nodes in $\text{path}(S_k)$. An example of eligibility is shown in Fig. 4.

We shall find t' in two phases. First we will efficiently locate a node $v \in \text{path}(S_k)$ of maximal depth such that there exists a relevant suffix S_t which is eligible at v . We will in fact prove that $S_{t'}$ is one of the suffixes that are eligible at v . Second, in the case that there are several relevant suffixes S_t which are eligible at v , we shall efficiently find $S_{t'}$ among them.

3.3.1. Finding Node v

Observe $\text{path}(S_k)$. We wish to search this path for the lowest node v for which there exists t such that S_t is relevant and eligible at v . Notice that the notion of eligibility satisfies the property that if some suffix S_t is eligible at some node $u \in \text{path}(S_k)$, then S_t is eligible at all of u 's ancestors as well. In other words, the eligibility property is monotone along $\text{path}(S_k)$. If the suffix tree had been preprocessed for answering level-ancestor queries, by the methods of, for example, [3], we can find the ancestor of ℓ_k of a specific depth d in $O(1)$ time. We conclude that we can perform a binary search on the depth of nodes on this path: at each node u we probe, we will efficiently test whether there exists some relevant suffix which is eligible at u .

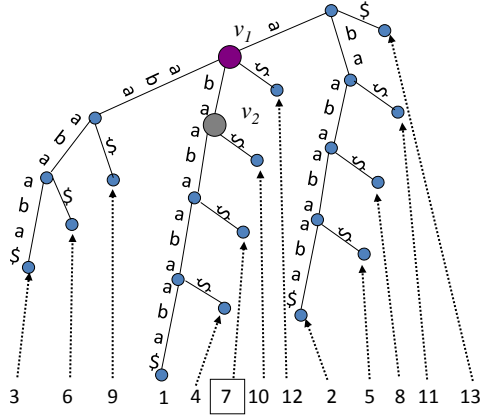


Figure 4: Again we look at the suffix tree of the string $S = abaabaabaaba\$$. This time we would like to encode the substring $S[7..9]$ within the context of $S[2..5]$. For this example, looking at location $r = 5$, the suffix S_4 is eligible at node v_1 and is not eligible at node v_2 .

We can perform this test by conducting the following range emptiness query: $\text{emptiness}([x(l_u), x(r_u)] \times [l, r - \text{length}(u) + 1])$. This is captured by the two following lemmas:

Lemma 1. *A suffix S_t is eligible at a node $u \in \text{path}(S_k)$ iff $u \in \text{path}(S_t)$ and $\text{length}(u) \leq r - t + 1$.*

PROOF. Let S^u be the string represented by $\text{path}(u)$. By definition, S_t is eligible at u if S^u is a prefix of $\text{LCP}_{r-t+1}(S_k, S_t)$. That happens if and only if (a) S^u is a prefix of $\text{LCP}(S_k, S_t)$ and therefore $u \in \text{path}(S_t)$, and (b) the length of S^u , i.e., $\text{length}(u)$ is at most $r - t + 1$.

Lemma 2. *Fix the range $[l, r]$ and let S_t be a suffix. Then $(x(\ell_t), y(\ell_t)) \in ([x(l_u), x(r_u)] \times [l, r - \text{length}(u) + 1])$ iff S_t is relevant and eligible at u .*

PROOF. Let S_t be a suffix. From the properties of a suffix tree, it holds that S_t 's lexicographical rank $x(\ell_t)$ is in $[x(l_u), x(r_u)]$ if and only if ℓ_t is a leaf in T_u or equivalently if and only if $u \in \text{path}(S_t)$. For its start location $t = y(\ell_t)$, $y(\ell_t) \in [l, r - \text{length}(u) + 1]$ if and only if both $y(\ell_t) \in [l, r]$ (i.e., S_t is relevant) and $y(\ell_t) \leq r - \text{length}(u) + 1$. As the latter can be re-stated as $\text{length}(u) \leq r - t + 1$, the lemma holds.

We conclude that the emptiness query returns a negative result if and only if there exists a suffix S_t which is both relevant and eligible at u .

Instead of the ordinary $O(\log n)$ -time binary search, we use a mixed “galloping” and ordinary binary search approach: we conduct the search by iterations, where in the i -th iteration we probe the node on the path whose depth is $2^{i-1} - 1$ and conduct the proper range emptiness query on it, repeating this process until

the first node whose emptiness query returned a positive result is encountered. Denote this node as q and denote the last node probed before q as p . Now we find v by binary searching on the sub-path between p and q . The main importance of the mixed search, is that now, using a refined analysis, we will later prove that encoding a phrase is done in time logarithmic in the phrase's length (rather than n).

Recall that we have defined $t' = \arg \max_{t \in [l, r]} |\text{LCP}_{r-t+1}(S_k, S_t)|$. At the end of this phase we have the following lemma:

Lemma 3. *v is the node of maximal depth at which $S_{t'}$ is eligible.*

PROOF. Assume not, and let $v' \neq v$ be the the node of maximal depth at which $S_{t'}$ is eligible. (Notice that v' is well-defined since $S_{t'}$ is trivially eligible at the root.) Since both $v, v' \in \text{path}(S_k)$ then either v is an ancestor of v' or vice-versa.

In the former case, $\text{depth}(v') > \text{depth}(v)$. However, by Lemma 2 and the monotonicity of the eligibility property, the binary search described before must conclude with the maximal depth node at which some relevant suffix is eligible, which contradicts the fact that $S_{t'}$ is eligible at v' .

In the latter case $\text{length}(v') < \text{length}(v)$ and by Lemma 2, there exists a relevant suffix S_t which is eligible at v , and therefore $|\text{LCP}_{r-t+1}(S_k, S_t)| \geq \text{length}(v)$. Since $\text{length}(v') < \text{length}(v)$ and by the definition of v' , $S_{t'}$ cannot be eligible at v . Therefore, $|\text{LCP}_{r-t'+1}(S_k, S_{t'})| < \text{length}(v)$, otherwise, when traversing $\text{LCP}_{r-t'+1}(S_k, S_{t'})$ starting from the root of the suffix tree, we would have to visit v , a fact which would make $S_{t'}$ eligible at v as well. Therefore $|\text{LCP}_{r-t+1}(S_k, S_t)| \geq \text{length}(v) > |\text{LCP}_{r-t'+1}(S_k, S_{t'})|$, which contradicts the optimality of t' .

3.3.2. Finding the Suffix $S_{t'}$

Recall that for a suffix S_t , $|\text{LCP}_{r-t+1}(S_k, S_t)| = \min\{|\text{LCP}(S_k, S_t)|, r-t+1\}$, and therefore this value can be computed in $O(1)$ using a single LCA query. Given v , as described in the previous phase, let $w \in \text{path}(S_k)$ be its specific child that is on $\text{path}(S_k)$. We then perform the following: we inspect v and query $\text{rangesucc}_y([x(l_v), x(r_v)] \times [l, \infty])$. Let (x', y') be the point returned by the query (see example in Fig. 5). We also inspect w and perform $\text{rangesucc}_y([x(l_w), x(r_w)] \times [l, \infty])$, and let (x'', y'') be the resulting point, if such exists. If (x'', y'') does not exist or that $y'' \notin [l, r - \text{length}(v) + 1]$, we choose t' to be y' . Otherwise, notice that y', y'' are the two start positions of the relevant suffixes $S_{y'}, S_{y''}$ respectively. Pick the value $t \in \{y', y''\}$ that maximizes the expression $|\text{LCP}_{r-t+1}(S_k, S_t)|$ and choose it to be t' (if both maximize the expression, choose $t' = \min\{y', y''\}$). Finally return the resulting substring, as represented by its start position t' and its length $|\text{LCP}_{r-t'+1}(S_k, S_{t'})|$. A pseudo-code of the above process can be seen in Procedure `Compute- t'` .

Before formally proving the correctness of the above procedure, we provide some intuition for it. Observe node v once again and the range $[x(l_v), x(r_v)] \times [l, r - \text{length}(v) + 1]$. There might be several relevant suffixes which are eligible at v , and therefore have their corresponding suffix points in that range. In this

Procedure Compute- t'

Input: v **Output:** selected t'

```
1  $w \leftarrow v$ 's child on path( $S_k$ );
2  $(x', y') \leftarrow \text{rangesucc}_y([x(l_v), x(r_v)] \times [l, \infty])$ ;
3  $(x'', y'') \leftarrow \text{rangesucc}_y([x(l_w), x(r_w)] \times [l, \infty])$ ;
4 if  $(x'', y'')$  does not exist or  $y'' \notin [l, r - \text{length}(v) + 1]$  then
5   |  $t' \leftarrow y'$ ;
6 else
7   | if  $|\text{LCP}_{r-y'+1}(S_k, S_{y'})| = |\text{LCP}_{r-y''+1}(S_k, S_{y''})|$  then
8     |  $t' \leftarrow \min\{y', y''\}$ ;
9     | else
10    |  $t' \leftarrow \arg \max_{t \in \{y', y''\}} |\text{LCP}_{r-t+1}(S_k, S_t)|$ ;
11 return  $t'$ ;
```

case, for each such suffix S_t , v is the node of maximal depth at which S_t is eligible. Notice that some of those suffixes may be represented by paths that branch out from v to w , and some by paths that branch out from v to one of its other children. Therefore, their potential “contribution to compression” may be different: for the specific case where for a suffix S_t , the path from the root to ℓ_t visits v and then w , and also it holds that $t + \text{length}(v) - 1 < r$, there is an additional eligible portion of S_t of length $r - (t + \text{length}(v) - 1)$ on the edge (v, w) (figuratively speaking; we mean of course that the additional eligible portion is a prefix of the substring represented by the label of (v, w)). We refer to such suffixes as *special suffixes* and note that the additional portion may be of a different length for different special suffixes. An example of the special suffix scenario is shown in Fig. 6. The existence of such special suffixes creates the need to inspect both v and w .

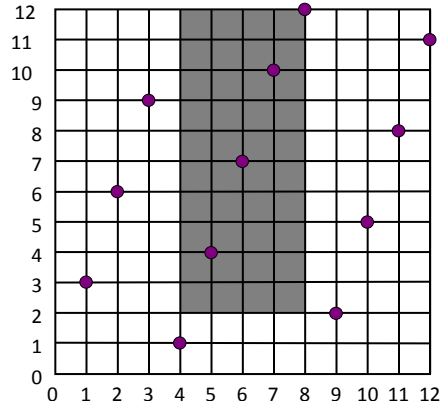


Figure 5: In Fig. 4 we found $v = v_1$, since v_1 is the node of maximal depth on the path from the root to S_7 , at which S_7 is eligible w.r.t $r = 5$. Therefore, the values of our range will be: $x(l_v) = 4$, $x(r_v) = 8$ and $l = 2$. Therefore, we obtain the range query as depicted by the grayed area. The output point of this query will be $(5, 4)$.

We move to the formal proof:

Theorem 1. *The value returned by the above procedure is $t' = \arg \max_{t \in [l, r]} |\text{LCP}_{r-t+1}(S_k, S_t)|$.*

PROOF. Consider all relevant suffixes that are eligible at v , and therefore are represented by suffix points in the range $[x(l_v), x(r_v)] \times [l, r - \text{length}(v) + 1]$ (by Lemma 2). By Lemma 3, $S_{t'}$ is among these suffixes, however we do not know whether $\text{path}(S_{t'})$ visits w , or not.

Notice that for all such suffixes S_t which branch out from v to a node other than w , it holds that $|\text{LCP}_{r-t+1}(S_k, S_t)| = \min\{|\text{LCP}(S_k, S_t)|, r - t + 1\} = |\text{LCP}(S_k, S_t)| = \text{length}(v)$. As these suffixes all have the same $|\text{LCP}_{r-t+1}(S_k, S_t)|$ value, our first rangesucc_y query simply picks the one with the leftmost start position.

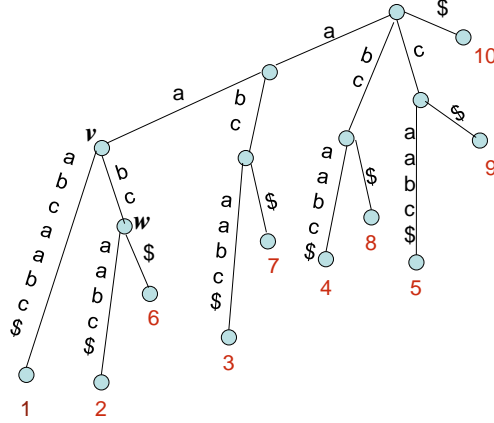


Figure 6: In this example, the suffix tree for the string $S = aaabcaabc\$$ is shown. Assume we wish to compress $S[6..9]$ (i.e., $k = 6$) and the given context is $S[1..4]$ (i.e., $r = 4$). Nodes v and w are labeled accordingly. Notice that the leftmost relevant suffix in T_v is $S_{y'} = S_1$ for which $|\text{LCP}_{r-y'+1}(S_k, S_{y'})| = \min\{\text{length}(v), r - y' + 1\} = 2$. However, the leftmost relevant suffix in T_w is $S_{y''} = S_2$ for which $|\text{LCP}_{r-y''+1}(S_k, S_{y''})| = r - y'' + 1 = 3$. Therefore we would like to output start location 2.

On the other hand, for all such suffixes S_t which branch out from v to w , it holds that $|\text{LCP}(S_k, S_t)| > r - t + 1$ (since for them, $|\text{LCP}(S_k, S_t)| \geq \text{length}(w)$, but they are not eligible at w) and therefore $|\text{LCP}_{r-t+1}(S_k, S_t)| = r - t + 1$. In other words, among these suffixes, the one maximizing $|\text{LCP}_{r-t+1}(S_k, S_t)|$ is the one with the leftmost start position. As all these suffixes have a one-to-one correspondence with suffix points in $[x(l_w), x(r_w)] \times [l, r - \text{length}(v) + 1]$ (by Lemma 2 and the fact that those suffixes have paths that visit w), our goal is to pick the suffix point in the range having the minimal y -coordinate. This is done by our second rangesucc_y query with the additional test that for the resulting point (x'', y'') (if such exists), indeed $y'' \in [l, r - \text{length}(v) + 1]$.

Since at first we do not know which of the cases holds for t' , the algorithm simply chooses the best of both results.

4. Outline of Substring Compression Query Algorithms

Given locations i and j which induce the substring $S[i..j]$ to be compressed, we describe the outline of our methods, in an inductive manner:

- For the first location i , two cases exist, according to query type:
 - SCQ:** write the encoded representation of $S[i]$.
 - GSCQ:** set $k \leftarrow i$ and calculate $\text{BLCP}(k, \alpha, \beta)$. For convenience, we denote $|\text{LCP}| = |\text{BLCP}(k, \alpha, \beta)|$.
- For a general location, assume $S[k..j]$ is left to be compressed. Again two cases exist:
 - SCQ:** the LZ method revolves around finding $\text{ILCP}(k, i, k - 1)$. For convenience, we denote $\text{LCP} = \text{ILCP}(k, i, k - 1)$.
 - GSCQ:** here we calculate both $\text{ILCP}(k, i, k - 1)$ and $\text{BLCP}(k, \alpha, \beta)$, and choose the longest of both. For convenience, this time we denote $|\text{LCP}| = \max\{|\text{ILCP}(k, i, k - 1)|, |\text{BLCP}(k, \alpha, \beta)|\}$.

It is important to note that in all cases we need not find the LCP itself, but rather it is sufficient to find its starting position t and its length. Once the proper $|\text{LCP}|$ value is obtained, if $k + |\text{LCP}| - 1 > j$, we truncate its last characters, leaving only the first $j - k + 1$.

If no such LCP exists (e.g., $|\text{ILCP}| = 0$ and, if applicable, $|\text{BLCP}| = 0$), we revert to writing the encoded representation of the current character, i.e., $S[k]$. Otherwise, we write the encoded representation of the distance to the starting position t (i.e., the value $k - t$) and length of LCP, and set $k \leftarrow k + |\text{LCP}|$. If $k \leq j$, we repeat this process, otherwise, we stop.

5. Substring Compression Query

Given a string $S[1..n]$, it will be preprocessed to efficiently answer queries of the form $\text{SCQ}(i, j)$, in which we are asked to find the compression of the substring $S[i..j]$. The compression of $S[i..j]$ will then be computed by performing ILCP queries in the manner described above until the compressed representation of the entire substring has been found.

5.1. Analysis

Our running times and space used are heavily affected by the choice of the range searching structure used. The following presents the general trade-off scheme:

Theorem 2. *$\text{SCQ}(i, j)$ can be answered in worst-case $O(C(i, j)Q_{\text{rsucc}})$ time, using a structure which employs $O(S_{\text{rsucc}})$ space, where Q_{rsucc} and S_{rsucc} stand for the query-time and space of the range successor structure, respectively.*

PROOF. We analyze the space and query time:

Table 2: SCQ tradeoffs

rsucc/rpred	Query Time	Space
[30]	$O(C(i, j) \log \log n)$	$O(n \log^\epsilon n)$
[30]	$O(C(i, j)(\log \log n)^2)$	$O(n \log \log n)$
[30]	$O(C(i, j) \log^\epsilon n)$	$O(n)$
[12]	$O(C(i, j))$	$O(n^{1+\epsilon})$

Space. Consists of: $O(n)$ for the suffix tree, augmented with the additional $x(\ell)$ and $\text{length}(u)$ values, and LCA information; S_{rsucc} for the range successor structure. Since in all of our configurations, the space for the range successor structure dominates the space requirements, we conclude the space used is $O(S_{\text{rsucc}})$.

Query Time. For each of the $C(i, j)$ phrases encoded, we use: Q_{rsucc} for range predecessor (resp. successor) queries made in order to find ℓ_{t_1} (resp. ℓ_{t_2}); $O(1)$ in order to compute both $|\text{LCA}(\ell_k, \ell_{t_1})|$ and $|\text{LCA}(\ell_k, \ell_{t_2})|$, and choose the maximum of both. We conclude the query time is overall $O(C(i, j)Q_{\text{rsucc}})$.

Substituting in the different flavors for range successor data structure of Nekrich and Navarro [30], we get the query time-space trade-offs presented in Table 2.

We also provide the following trade-off:

Theorem 3. *For any constant $\epsilon > 0$, $SCQ(i, j)$ can be answered in worst-case $O(C(i, j))$ time, using a structure which employs $O(n^{1+\epsilon})$ space.*

PROOF. Notice that our range queries are performed on $x(\ell)$ and $y(\ell)$ values. A unique property of these values is that no $x(\ell)$ or $y(\ell)$ value occurs twice in P , i.e., the sequence of point x -coordinates, and the sequence of point y -coordinates, are both permutations of $[n + 1]$. Using the *range next value* structure of [12] allows us to obtain the following tradeoff: the space used is dominated by the $O(n^{1+\epsilon})$ space required for the range successor structure, and for the query time, since a single range successor/predecessor query can now be answered in $O(1)$, the overall query time is worst-case $O(C(i, j))$.

6. General Substring Compression Query

For GSCQ, in addition to the two locations i and j , which denote the substring $S[i..j]$ to be compressed, we receive two more indices α and β , which induce a *context* substring $S[\alpha.. \beta]$. This time we are asked to provide $\text{LZ}(S[i..j] \mid S[\alpha.. \beta])$.

Here, when trying to compress $S[k..j]$ for some $i \leq k \leq j$, we have two options: for the first we consider phrases having a start position $i \leq t \leq k - 1$. This option is the one solved in Section 5, using ILCP queries. The second, is to consider phrases taken from $S[\alpha.. \beta]$. This will be done using a BLCQ query.

6.1. Analysis

The analysis is depicted in the following theorem:

Theorem 4. $\text{GSCQ}(i, j, \alpha, \beta)$ can be answered in worst-case

$$O(C_{\alpha, \beta}(i, j) \left(\log\left(\frac{j-i}{C_{\alpha, \beta}(i, j)}\right) Q_{\text{empt}} + Q_{\text{rsucc}} \right))$$

time, using a structure which takes $O(S_{\text{empt}} + S_{\text{rsucc}})$ space, where Q_{empt} and S_{empt} (resp. Q_{rsucc} and S_{rsucc}) are the query time and space of the range emptiness (resp. range successor) structure, respectively.

PROOF. As follows:

Space. Consists of: $O(n)$ for the suffix tree, augmented with the additional $x(l_u)$, $x(r_u)$ and $\text{length}(u)$ values, LCA and level-ancestor structure information. These bounds will be dominated by the range searching structures chosen.

Query Time. Consider the query's main loop described in Section 4 and consider the d -th iteration of the query algorithm main loop, and let len_d be the length of the phrase encoded in this iteration ($d = 1, \dots, C_{\alpha, \beta}(i, j)$). Assume $S[k..j]$ is the portion left to be compressed just before this iteration, and let nodes v , p , and q be as defined before. It holds that $\text{depth}(p) \leq \text{length}(p) \leq \text{length}(v) \leq |\text{BLCP}(k, \alpha, \beta)|$. Node q was found one iteration after node p . Therefore:

$$\text{depth}(q) \leq 2(\text{depth}(p) + 1) \leq 2(|\text{BLCP}(k, \alpha, \beta)| + 1). \quad (1)$$

We conclude that finding q was done by performing $O(\log|\text{BLCP}(k, \alpha, \beta)|)$ node accesses, and the following binary search, was supported by performing

$$O(\log(\text{depth}(q) - \text{depth}(p))) = O(\log|\text{BLCP}(k, \alpha, \beta)|) \quad (2)$$

node accesses. Since

$$|\text{BLCP}(k, \alpha, \beta)| \leq \max\{|\text{ILCP}(k, i, k-1)|, |\text{BLCP}(k, \alpha, \beta)|\} = \text{len}_d, \quad (3)$$

and when accessing each node, a range emptiness query was conducted, overall time for the mixed search described is $O(\log(\text{len}_d) \cdot Q_{\text{empt}})$, where Q_{empt} is the query time used for the emptiness query. We conclude that a $\text{BLCP}(k, \alpha, \beta)$ query can be answered in $O(\log(\text{len}_d) \cdot Q_{\text{empt}} + Q_{\text{rsucc}})$, where Q_{rsucc} is the time required for each of the two final range successor queries performed. Recall that an $\text{ILCP}(k, i, k-1)$ query is also made, however, this query only takes $O(Q_{\text{empt}})$ time, and therefore does not influence the time bound.

We conclude that GSCQ can be answered in overall

$$O\left(Q_{\text{empt}} \sum_{d=1}^{C_{\alpha, \beta}(i, j)} \log(\text{len}_d) + C_{\alpha, \beta}(i, j) Q_{\text{rsucc}}\right) \quad (4)$$

Table 3: GSCQ tradeoffs

empt	rsucc	Query Time	Space
[12]	[12]	$O(C_{\alpha,\beta}(i,j) \log(\frac{j-i}{C_{\alpha,\beta}(i,j)}))$	$O(n^{1+\epsilon})$
[5]	[30]	$O(C_{\alpha,\beta}(i,j) \log(\frac{j-i}{C_{\alpha,\beta}(i,j)}) \log \log n)$	$O(n \log^\epsilon n)$
[5]	[30]	$O(C_{\alpha,\beta}(i,j) (\log(\frac{j-i}{C_{\alpha,\beta}(i,j)}) \log \log n + (\log \log n)^2))$	$O(n \log \log n)$
[30]	[30]	$O(C_{\alpha,\beta}(i,j) \log(\frac{j-i}{C_{\alpha,\beta}(i,j)}) \log^\epsilon n)$	$O(n)$

time. $\{\text{len}_d\}_{d=1}^{C_{\alpha,\beta}(i,j)}$ is a partition of $|S[i..j]| = j - i + 1$, therefore the above expression is maximized when $\text{len}_1 = \dots = \text{len}_{C_{\alpha,\beta}(i,j)} = \frac{j-i+1}{C_{\alpha,\beta}(i,j)}$. We conclude that $\text{GSCQ}(i, j, \alpha, \beta)$ can be answered in

$$O(C_{\alpha,\beta}(i,j) (\log(\frac{j-i}{C_{\alpha,\beta}(i,j)}) Q_{\text{empt}} + Q_{\text{rsucc}})) \quad (5)$$

time.

The choice of range emptiness and range successor structures will determine the time bounds for their respective queries. Tradeoff results are given in the following corollaries, and in Table 3, where the column labeled “empt” denotes the range emptiness structure used, and the column labeled “rsucc” denotes the range successor structure used.

7. Conclusion

The goal of our research was to find the inherent connection between the compressed representation of a string and that of its substrings. We have based our research on the Lempel-Ziv algorithm and focused our attention on two main variants of the problem: basic substring compression, and generalized substring compression. For these problems we have achieved the following query times: given a string S , for a chosen substring $S[i..j]$ (possibly with the substring $S[\alpha.. \beta]$ as a context), assuming C phrases are needed for the compressed representation of $S[i..j]$ using the LZ77 algorithm, the best query times we achieve are $O(C)$ and $O(C \log(\frac{j-i}{C}))$ for the basic substring compression query and the generalized substring compression query, respectively. While initial results for these problems were presented in [8], our results are an improvement for the basic substring compression query, and are the first known correct ones for the generalized substring compression query. The problems we have dealt with, and the solutions proposed, leave us with two main problems that may each be of independent interest.

1. The main problem left open following our research pertains to other compression algorithms. Specifically, how can other compression mechanisms

be adapted to deal with substring compression. This problem was presented in the original paper dealing with this problem [8], and we leave it open following our research as well. It seems as though algorithms such as Arithmetic encoding and algorithms relying on the Burrows-Wheeler transform would be more difficult to adapt to substring compression if at all possible, however this is simply a conjecture and such research should be of independent interest.

2. Another interesting problem may be to investigate the solutions for substrings for other classic stringology problems. For example, in the classical edit-distance problem, how can we efficiently retrieve the edit distance for any two substrings of the strings at hand. Specifically, can the known solutions inherently solve the problem for substrings as well? Or how can they be adapted to work for substrings?

References

- [1] S. Alstrup, G. S. Brodal, and T. Rauhe. New data structures for orthogonal range searching. In *IEEE Symposium on Foundations of Computer Science*, pages 198–207, 2000.
- [2] A. Amir, D. Keselman, G. M. Landau, M. Lewenstein, N. Lewenstein, and M. Rodeh. Text indexing and dictionary matching with one error. *J. Algorithms*, 37(2):309–325, 2000.
- [3] M. A. Bender and M. Farach-Colton. The level ancestor problem simplified. *Theor. Comput. Sci.*, 321(1):5–12, 2004.
- [4] P. Bille and I. L. Gørtz. Substring range reporting. In *Proc. of Symposium on Combinatorial Pattern Matching (CPM)*, pages 299–308, 2011.
- [5] T. M. Chan, K. G. Larsen, and M. Patrascu. Orthogonal range searching on the ram, revisited. In *Symposium on Computational Geometry*, pages 1–10, 2011.
- [6] M. Charikar, E. Lehman, D. Liu, R. Panigrahy, M. Prabhakaran, A. Rasala, A. Sahai, and A. Shelat. Approximating the smallest grammar: Kolmogorov complexity in natural models. In J. H. Reif, editor, *STOC*, pages 792–801. ACM, 2002.
- [7] X. Chen, S. Kwong, and M. Li. A compression algorithm for dna sequences and its applications in genome comparison. In *RECOMB*, page 107, 2000.
- [8] G. Cormode and S. Muthukrishnan. Substring compression problems. In *SODA '05: Proceedings of the sixteenth annual ACM-SIAM Symposium on Discrete Algorithms*, pages 321–330, Philadelphia, PA, USA, 2005. Society for Industrial and Applied Mathematics.

- [9] A. J. Cox, T. Jakobi, G. Rosone, and O. Schulz-Trieglaff. Comparing dna sequence collections by direct comparison of compressed text indexes. In B. J. Raphael and J. Tang, editors, *WABI*, volume 7534 of *Lecture Notes in Computer Science*, pages 214–224. Springer, 2012.
- [10] M. Crochemore and L. Ilie. Computing longest previous factor in linear time and applications. *Inf. Process. Lett.*, 106(2):75–80, 2008.
- [11] M. Crochemore, L. Ilie, and W. F. Smyth. A simple algorithm for computing the lempel ziv factorization. In *DCC*, pages 482–488. IEEE Computer Society, 2008.
- [12] M. Crochemore, C. S. Iliopoulos, M. Kubica, M. S. Rahman, and T. Walen. Improved algorithms for the range next value problem and applications. In *STACS 2008, 25th Annual Symposium on Theoretical Aspects of Computer Science*, volume 1 of *LIPICs*, pages 205–216. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, Germany, 2008.
- [13] M. Farach. Optimal suffix tree construction with large alphabets. In *FOCS '97: Proceedings of the 38th Annual Symposium on Foundations of Computer Science (FOCS '97)*, page 137, Washington, DC, USA, 1997. IEEE Computer Society.
- [14] P. Ferragina. Dynamic text indexing under string updates. *J. Algorithms*, 22(2):296–328, 1997.
- [15] P. Ferragina, S. Muthukrishnan, and M. de Berg. Multi-method dispatching: A geometric approach with applications to string matching problems. In *STOC*, pages 483–491, 1999.
- [16] K. Goto and H. Bannai. Simpler and faster lempel ziv factorization. *CoRR*, abs/1211.3642, 2012.
- [17] D. Harel and R. E. Tarjan. Fast algorithms for finding nearest common ancestors. *SIAM J. Computing*, 13(2):338–355, 1984.
- [18] D. Hermelin, S. Landau, G. Landau, and O. Weimann. A unified algorithm for accelerating edit-distance via text-compression. In *Proc. of the 26th International Symposium on Theoretical Aspects of Computer Science (STACS)*, pages 529–540, 2009.
- [19] W.-K. Hon, R. Shah, S. V. Thankachan, and J. S. Vitter. On position restricted substring searching in succinct space. *Journal of Discrete Algorithms*, 17:109–114, 2012.
- [20] J. Kärkkäinen, D. Kempa, and S. J. Puglisi. Linear time lempel-ziv factorization: Simple, fast, small. In J. Fischer and P. Sanders, editors, *CPM*, volume 7922 of *Lecture Notes in Computer Science*, pages 189–200. Springer, 2013.

- [21] O. Keller, T. Kopelowitz, S. Landau, and M. Lewenstein. Generalized substring compression. In *Combinatorial Pattern Matching, 20th Annual Symposium, CPM 2009*, volume 5577 of *Lecture Notes in Computer Science*, pages 26–38. Springer, 2009.
- [22] O. Keller, T. Kopelowitz, and M. Lewenstein. Range non-overlapping indexing and successive list indexing. In *WADS '07: Proceedings of the 10th International Workshop on Algorithms and Data Structures*, volume 4619 of *Lecture Notes in Computer Science*, pages 626–631. Springer-Verlag, 2007.
- [23] T. Kopelowitz, M. Lewenstein, and E. Porat. Persistency in suffix trees with applications to string interval problems. In *Proc. of Symposium on String Processing and Information Retrieval (SPIRE)*, pages 67–80, 2011.
- [24] S. Kuruppu, S. J. Puglisi, and J. Zobel. Relative lempel-ziv compression of genomes for large-scale storage and retrieval. In E. Chávez and S. Lonardi, editors, *SPIRE*, volume 6393 of *Lecture Notes in Computer Science*, pages 201–206. Springer, 2010.
- [25] H.-P. Lenhof and M. Smid. Using persistent data structures for adding range restrictions to searching problems. *RAIRO Theoretical Informatics and Applications*, 28:25–49, 1994.
- [26] M. Lewenstein. Orthogonal range searching for text indexing. *CoRR*, abs/1306.0615, 2013.
- [27] V. Mäkinen and G. Navarro. Position-restricted substring searching. In *LATIN 2006: Theoretical Informatics, 7th Latin American Symposium*, volume 3887 of *Lecture Notes in Computer Science*, pages 703–714. Springer, 2006.
- [28] E. M. McCreight. A space-economical suffix tree construction algorithm. *J. ACM*, 23(2):262–272, 1976.
- [29] S. Muthukrishnan. Personal communication with the second author, 2005.
- [30] Y. Nekrich and G. Navarro. Sorted range reporting. In F. V. Fomin and P. Kaski, editors, *SWAT*, volume 7357 of *Lecture Notes in Computer Science*, pages 271–282. Springer, 2012.
- [31] E. Rivals, O. Delgrange, J.-P. Delahaye, M. Dauchet, M.-O. Delorme, A. Hénaut, and E. Ollivier. Detection of significant patterns by compression algorithms: the case of approximate tandem repeats in dna sequences. *Computer Applications in the Biosciences*, 13(2):131–136, 1997.
- [32] W. Rytter. Application of Lempel-Ziv factorization to the approximation of grammar-based compression. *Theoretical Computer Science*, 302(1-3):211–222, 2003.

- [33] E. Ukkonen. On-line construction of suffix trees. *Algorithmica*, 14(3):249–260, 1995.
- [34] P. Weiner. Linear pattern matching algorithms. In *14th Annual Symposium on Switching and Automata Theory*, pages 1–11. IEEE, 1973.
- [35] J. Ziv and A. Lempel. A universal algorithm for sequential data compression. *IEEE Transactions on Information Theory*, 23(3):337–343, 1977.