# Closed Factorization[☆]

Golnaz Badkobeh[1], Hideo Bannai[2], Keisuke Goto[2], Tomohiro I[3],
Costas S. Iliopoulos[4], Shunsuke Inenaga[2], Simon J. Puglisi[5], Shiho Sugimoto[2]

## Abstract

A *closed string* is a string with a proper substring that occurs in the string as a
prefix *and* a suffix, but not elsewhere. Closed strings were introduced by Fici (Proc.
WORDS, 2011) as objects of combinatorial interest in the study of Trapezoidal and
Sturmian words. In this paper we present algorithms for computing closed factors
(substrings) in strings. First, we consider the problem of greedily factorizing a string
into a sequence of longest closed factors. We describe an algorithm for this problem
that uses linear time and space. We then consider the related problem of computing,
for every position in the string, the longest closed factor starting at that position.
We describe a simple algorithm for the problem that runs in $O(n \log n / \log \log n)$
time, where $n$ is the length of the string. This also leads to an algorithm to compute
the maximal closed factor containing (i.e. covering) each position in the string in
$O(n \log n / \log \log n)$ time. We also present linear time algorithms to factorize a
string into a sequence of shortest closed factors of length at least two, to compute
the shortest closed factor of length at least two starting at each position of the
string, and to compute a minimal closed factor of length at least two containing
each position of the string.

*Keywords:*   closed strings, borders, suffix trees, range successor queries,
Manhattan skyline problem

## 1. Introduction

A *closed string* is a string with a proper substring that occurs as a prefix and a
suffix but not elsewhere in the string. Closed strings were introduced by Fici [1] as
objects of combinatorial interest in the study of Trapezoidal and Sturmian words.
Since then, Badkobeh, Fici, and Liptak [2] have proved a tight lowerbound for the
number of closed factors (substrings) in strings of given length and alphabet size.

In this paper we initiate the study of algorithms for computing closed factors.
In particular we consider the following algorithmic problems.

[1]Department of Computer Science, University of Sheffield, United Kingdom
[2]Department of Informatics, Kyushu University, Japan
[3]Department of Computer Science, TU Dortmund, Germany
[4]Department of Informatics, King's College London, United Kingdom
[5]Helsinki Institute for Information Technology, Department of Computer Science, University
of Helsinki, Finland

**Longest closed factorization problem** This problem is to factorize a given string into a sequence of longest closed factors (we give a formal definition of the problem below, in Section 2). We describe an algorithm for this problem that uses $O(n)$ time and space, where $n$ is the length of the given string.

**Longest closed factor array problem** This problem requires us to compute the length of the longest closed factor starting at each position in the input string. We show that this problem can be solved in $O(n\frac{\log n}{\log \log n})$ time with $O(n)$ space, using techniques from computational geometry.

**Maximal closed factor array problem** This problem asks us to compute the length of a longest closed factor containing (i.e. covering) each position in the input string. We show that this problem can also be solved in $O(n\frac{\log n}{\log \log n})$ time with $O(n)$ space.

**Shortest closed factorization problem** This is the problem of factorizing a given string into a sequence of shortest closed factors of length at least two. We show this can be solved in $O(n)$ time and $O(\sigma)$ space, where $\sigma$ is the alphabet size.

**Shortest closed factor array problem** This problem requires us to compute the length of the shortest closed factor of length at least two starting at each position in the input string. We solve this problem in $O(n)$ time and space.

**Minimal closed factor array problem** This problem asks us to compute, for each position in the input string, the length of the shortest closed factor of length at least two containing the position. We present an $O(n)$-time and space algorithm for this problem.

This paper proceeds as follows. In the next section we set notation, define the problems more formally, and outline basic data structures and concepts. Section 3 describes an efficient solution to the longest closed factorization problem and Section 4 then considers the longest closed factor array and the maximal closed factor array. In Section 5, we present efficient algorithms for computing the shortest closed factorization, the shortest closed factor array, and the minimal closed factor array. Reflections and outlook are offered in Section 6.

Some of these results appeared in a preliminary version of this paper [3].

## 2. Preliminaries

### 2.1. Strings and Closed Factorization

Let $\Sigma$ denote a fixed alphabet of $|\Sigma|$ distinct letters (or characters). An element of $\Sigma^*$ is called a string. For any strings $W, X, Y, Z$ such that $W = XYZ$, the strings $X, Y, Z$ are respectively called a prefix, substring, and suffix of $W$. The length of a string $X$ will be denoted by $|X|$. Let $\varepsilon$ denote the empty string of length 0, i.e., $|\varepsilon| = 0$. For any non-negative integer $n$, $X[1, n]$ denotes a string $X$ of length $n$. A prefix $X$ of a string $W$ with $|X| < |W|$ is called a proper prefix of $W$. Similarly, a suffix $Z$ of $W$ with $|Z| < |W|$ is called a proper suffix of $W$. For any string $X$ and integer $1 \leq i \leq |X|$, let $X[i]$ denote the $i$th character of $X$, and for any integers $1 \leq i \leq j \leq |X|$, let $X[i..j]$ denote the substring of $X$ that starts at position $i$ and ends at position $j$. For convenience, let $X[i..j]$ be the empty string if $j < i$. For any strings $X$ and $Y$, if $Y = X[i..j]$, then we say that $i$ is an occurrence of $Y$ in $X$.

If a non-empty string $X$ is both a proper prefix and suffix of string $W$, then, $X$ is called a *border* of $W$. A string $W$ is said to be *closed*, if there exists a border $X$ of $W$ that occurs exactly twice in $W$, i.e., $X = W[1..|X|] = W[|W| - |X| + 1..|W|]$ and

$X \neq W[i..i + |X| - 1]$ for any $2 \leq i \leq |W| - |X|$. We define a single character $C \in \Sigma$ to be closed, assuming that the empty string $\varepsilon$ occurs exactly twice in $C$. A string $X$ is a *closed factor* of $W$, if $X$ is closed and is a substring of $W$. Throughout we consider a string $X[1..n]$ on $\Sigma$. We define the *longest closed factorization* of string $X[1..n]$ as follows.

**Definition 1 (Longest Closed Factorization)** *The longest closed factorization of string $X[1..n]$, denoted $\mathsf{LCF}(X)$, is a sequence $(\mathsf{G}_0, \mathsf{G}_1, \ldots, \mathsf{G}_k)$ of strings such that $\mathsf{G}_0 = \varepsilon$, $X[1, n] = \mathsf{G}_0\mathsf{G}_1 \cdots \mathsf{G}_k$ and, for each $1 \leq j \leq k$, $\mathsf{G}_j$ is the longest prefix of $X[|\mathsf{G}_0 \cdots \mathsf{G}_{j-1}| + 1..n]$ that is closed.*

**Example 1** *For string* $X = \mathtt{ababaacbbbcbcc}$, $\mathsf{LCF}(X) = (\varepsilon, \mathtt{ababa}, \mathtt{a}, \mathtt{cbbbcb}, \mathtt{cc})$.

We remark that a closed factor $\mathsf{G}_j$ in an LCF is a single character if and only if $|\mathsf{G}_1 \cdots \mathsf{G}_{j-1}| + 1$ is the rightmost (last) occurrence of character $X[|\mathsf{G}_1 \cdots \mathsf{G}_{j-1}| + 1]$ in $X$.

We also define the *longest closed factor array* of string $X[1..n]$.

**Definition 2 (Longest Closed Factor Array)** *The longest closed factor array of string $X[1..n]$ is an array $\mathsf{LNG}[1..n]$ of integers such that for any $1 \leq i \leq n$, $\mathsf{LNG}[i] = \ell$ if and only if $\ell$ is the length of the longest prefix of $X[i..n]$ that is closed.*

**Example 2** *For string* $X = \mathtt{ababaacbbbcbcc}$, $\mathsf{LNG} = [5, 4, 3, 5, 2, 1, 6, 3, 2, 4, 3, 1, 2, 1]$.

Clearly, given the longest closed factor array $\mathsf{LNG}[1..n]$ of string $X$, $\mathsf{LCF}(X)$ can be computed in $O(n)$ time. However, the algorithm we describe in Section 4 to compute $\mathsf{LNG}[1..n]$ requires $O(n\frac{\log n}{\log \log n})$ time, and so using it to compute $\mathsf{LCF}(X)$ would also take $O(n\frac{\log n}{\log \log n})$ time overall. In Section 3 we present an optimal $O(n)$-time algorithm to compute $\mathsf{LCF}(X)$ that does not require $\mathsf{LNG}[1..n]$.

**Definition 3 (Maximal Closed Factor Array)** *The maximal closed factor array of string $X[1..n]$ is an array $\mathsf{MAX}[1..n]$ of integers such that for any $1 \leq i \leq n$, $\mathsf{MAX}[i] = \ell$ if and only if $\ell$ is the length of a longest closed factor of $X$ that contains position $i$.*

**Example 3** *For string* $X = \mathtt{ababaacbbbcbcc}$, $\mathsf{MAX} = [5, 5, 5, 5, 5, 5, 6, 6, 6, 6, 6, 6, 3, 2]$.

As counterparts of the longest closed factorization and the longest closed factor array, we respectively define the *shortest closed factorization* and the *shortest closed factor array* of string $X[1..n]$ below. If single characters are allowed to be used as closed factors, then the problems of computing the shortest closed factorization and the shortest closed factor array become trivial. To make the problems more interesting, we only consider closed factors of length at least 2.

**Definition 4 (Shortest Closed Factorization)** *The shortest closed factorization of string $X[1..n]$, denoted $\mathsf{SCF}(X)$, is a sequence $(\mathsf{F}_0, \mathsf{F}_1, \ldots, \mathsf{F}_h)$ of strings such that $\mathsf{F}_0 = \varepsilon$, $X[1..n] = \mathsf{F}_0\mathsf{F}_1 \cdots \mathsf{F}_h$ and, for each $1 \leq j \leq h$, $\mathsf{F}_j$ is the shortest prefix of $X[|\mathsf{F}_0 \cdots \mathsf{F}_{j-1}| + 1..n]$ that is closed and is of length at least 2.*

**Example 4** *For string* $X = \mathtt{ababaacbbbcbcc}$, $\mathsf{SCF}(X) = (\varepsilon, \mathtt{aba}, \mathtt{baacb}, \mathtt{bb}, \mathtt{bcb}, \mathtt{cc})$.
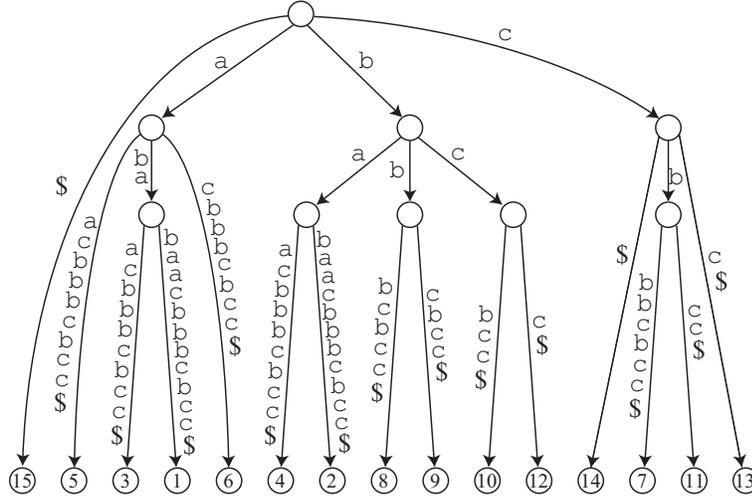
Figure 1: The suffix tree of string $X = \texttt{ababaacbbbcbcc\$}$, where each leaf stores the beginning position of the corresponding suffix. All the branches from an internal node are sorted in ascending lexicographical order, assuming $ is the lexicographically smallest. The suffix array SA of X is $[15, 5, 3, 1, 6, 4, 2, 8, 9, 10, 12, 14, 7, 11, 13]$, which corresponds to the sequence of leaves from left to right, and thus can be computed in linear time by a depth first traversal on the suffix tree.

**Definition 5 (Shortest Closed Factor Array)** *The shortest closed factor array of string $X[1..n]$ is an array $\mathsf{SHT}[1..n]$ of integers such that for any $1 \leq i \leq n$, $\mathsf{SHT}[i] = \ell$ if $\ell$ is the length of the shortest prefix of $X[i..n]$ that is closed and is of length at least 2, and $\mathsf{SHT}[i] = null$ if there are no prefixes of $X[i..n]$ that are closed and are of length at least 2.*

**Example 5** *For string* $X = \texttt{ababaacbbbcbcc}$, $\mathsf{SHT} = [3, 3, 3, 5, 2, null, 5, 2, 2, 3, 3, null, 2, null]$.

**Definition 6 (Minimal Closed Factor Array)** *The minimal closed factor array of string $X[1..n]$ is an array $\mathsf{MIN}[1..n]$ of integers such that for any $1 \leq i \leq n$, $\mathsf{MIN}[i] = \ell$ if $\ell$ is the length of a shortest closed factor of X of length at least 2 that contains position $i$ if it exists, and $\mathsf{MIN}[i] = null$ otherwise.*

**Example 6** *For string* $X = \texttt{ababaacbbbcbcc}$, $\mathsf{MIN} = [3, 3, 3, 3, 3, 5, 5, 2, 2, 2, 3, 3, 2, 2]$.

*2.2. Tools*

The suffix array [4] $\mathsf{SA_X}$ (we drop subscripts when they are clear from the context) of a string X is an array $\mathsf{SA}[1..n]$ which contains a permutation of the integers $[1..n]$ such that $X[\mathsf{SA}[1]..n] < X[\mathsf{SA}[2]..n] < \cdots < X[\mathsf{SA}[n]..n]$. In other words, $\mathsf{SA}[j] = i$ iff $X[i..n]$ is the $j^{\text{th}}$ suffix of X in ascending lexicographical order.

The suffix tree [5] of a string $X[1..n]$ is a compacted trie consisting of all suffixes of X. Suffix trees can be represented in linear space, and can be constructed in linear time for integer alphabets [6]. Figure 1 illustrates the suffix tree for an example string.

For a node $w$ in the suffix tree let $\mathsf{pathlabel}(w)$ be the string spelt out by the letters on the edges on the path from the root to $w$. If there is a branch from node $u$ to node $v$ and $u$ is an ancestor of $v$ then we say $u = \mathsf{parent}(v)$. Assuming that every string X terminates with a special character $ which occurs nowhere else in X, there is a one-to-one correspondence between the suffixes of X and the leaves of the suffix tree of X. We assume the branches from a node $u$ to each child $v$ of $u$ are stored in ascending lexicographical order of $\mathsf{pathlabel}(v)$. When this is the case,

4

SA is simply the leaves of the suffix tree when read during a depth-first traversal. At each internal node $v$ in the suffix tree we store two additional values $v.s$ and $v.e$ such that $\mathsf{SA}[v.s..v.e]$ contains the beginning positions of all the suffixes in the subtree rooted at $v$.

## 3. Greedy Longest Closed Factorization in Linear Time

In this section, we present an algorithm to compute the closed factorization $\mathsf{LCF}(\mathsf{X})$ of a given string $\mathsf{X}[1..n]$. Our high level strategy is to build a data structure that helps us to efficiently compute, for a given position $i$ in $\mathsf{X}$, the longest closed factor starting at $i$. The core of this data structure is the suffix tree for $\mathsf{X}$, which we decorate in various ways.

Let $S$ be the set of the beginning positions of the longest closed factors in $\mathsf{LCF}(\mathsf{X})$. For any $i \in S$, let $\mathsf{G} = \mathsf{X}[i..i + |\mathsf{G}| - 1]$ be the longest closed factor of $\mathsf{X}$ starting at position $i$ in $\mathsf{X}$.

Let $\mathsf{G}'$ be the unique border of the longest closed factor $\mathsf{G}$ starting at position $i$ of $\mathsf{X}$, and $b_i$ be its length, i.e., $\mathsf{G}' = \mathsf{G}[1..b_i] = \mathsf{X}[i..i+b_i-1]$ (if $\mathsf{G}$ is a single character, then $\mathsf{G}' = \varepsilon$ and $b_i = 0$). The following lemma shows that we can efficiently compute $\mathsf{LCF}(\mathsf{X})$ if we know $b_i$ for all $i \in S$.

**Lemma 1** *Given $b_i$ for all $i \in S$, we can compute $\mathsf{LCF}(\mathsf{X})$ in a total of $\mathrm{O}(n)$ time and space independently of the alphabet size.*

PROOF. If $b_i = 0$, then $\mathsf{G} = \mathsf{X}[i]$. Hence, in this case it clearly takes $\mathrm{O}(1)$ time and space to compute $\mathsf{G}$.

If $b_i \geq 1$, then we can compute $\mathsf{G}$ in $\mathrm{O}(|\mathsf{G}|)$ time and $\mathrm{O}(b_i)$ space, as follows. We preprocess the border $\mathsf{G}'$ of $\mathsf{G}$ using the Knuth-Morris-Pratt (KMP) string matching algorithm [7]. This preprocessing takes $\mathrm{O}(b_i)$ time and space. We then search for the first occurrence of $\mathsf{G}'$ in $\mathsf{X}[i + 1..n]$ (i.e. the next occurrence of the longest border of $\mathsf{G}[1..m]$ to the right of the occurrence $\mathsf{X}[i..i + b_i - 1]$). The location of the next occurrence tells us where the end of the closed factor is, and so it also tells us $\mathsf{G} = \mathsf{X}[i..i + |\mathsf{G}| - 1]$. The search takes $\mathrm{O}(|\mathsf{G}|)$ time — i.e. time proportional to the length of the closed factor. Because the sum of the lengths of the closed factors is $n$, over the whole factorization we take $\mathrm{O}(n)$ time and space. The running time and space usage of the algorithm are clearly independent of the alphabet size. □

What remains is to be able to efficiently compute $b_i$ for a given $i \in S$. The following lemma gives an efficient solution to this subproblem.

**Lemma 2** *We can preprocess the suffix tree of string $\mathsf{X}[1, n]$ in $\mathrm{O}(n)$ time and space, so that $b_i$ for each $i \in S$ can be computed in $\mathrm{O}(1)$ time.*

PROOF. In each leaf of the suffix tree, we store the beginning position of the suffix corresponding to the leaf. For any internal node $v$ of the suffix tree of $\mathsf{X}$, let $\max(v)$ denote the maximum leaf value in the subtree rooted at $v$, i.e.,

$$\max(v) = \max\{i \mid \mathsf{X}[i..i+|\mathsf{pathlabel}(v)|-1] = |\mathsf{pathlabel}(v)|, 1 \leq i \leq n-|\mathsf{pathlabel}(v)|-1\}.$$

We can compute $\max(v)$ for every $v$ in a total of $\mathrm{O}(n)$ time total via a depth first traversal. Next, let $\mathsf{P}[1..n]$ be an array of pointers to suffix tree nodes (to be computed next). Initially every $\mathsf{P}[i]$ is set to null. We traverse the suffix tree in pre-order, and for each node $v$ we encounter we set $\mathsf{P}[\max(v)] = v$ if $\mathsf{P}[\max(v)]$ is null. At the end of the traversal $\mathsf{P}[i]$ will contain a pointer to the highest node $w$ in

Figure 2: Illustration for Lemma 2. We consider the longest closed factor $\mathsf{G}$ starting at position $i$ of string $\mathsf{X}$. We retrieve node $\mathsf{P}[i] = v$, which implies $\max(v) = i$. Let $u$ be the parent of $v$. The black circle represents a (possibly implicit) node which represents $\mathsf{X}[i..i + |\mathsf{pathlabel}(u)| - 1]$, which has the same set of occurrences as $\mathsf{pathlabel}(v)$. Hence $b_i = |\mathsf{pathlabel}(u)|$, and therefore $\mathsf{G} = \mathsf{X}[i..j + |\mathsf{pathlabel}(u)| - 1]$, where $j$ is the leftmost occurrence of $\mathsf{pathlabel}(u)$ with $j > i$.

the tree for which $i$ is the maximum leaf value (i.e., $i$ is the rightmost occurrence of $\mathsf{pathlabel}(w)$).

We are now able to compute $b_i$, the length of the unique border of the longest closed factor starting at any given $i$, as follows. First we retrieve node $v = \mathsf{P}[i]$. Observe that, because of the definition of $\mathsf{P}[i]$, there are no occurrences of substring $\mathsf{X}[i..i + |\mathsf{pathlabel}(v)| - 1]$ to the right of $i$. Let $u = \mathsf{parent}(v)$ (note $v$ cannot be the root because it has non-empty $\mathsf{pathlabel}$). There are two cases to consider:

- If $u$ is not the root, then observe that there always exists an occurrence of substring $\mathsf{pathlabel}(u)$ to the right of position $i$ (otherwise $i$ would be the rightmost occurrence of $\mathsf{pathlabel}(u)$, but this cannot be the case since $u$ is higher than $v$, and we defined $\mathsf{P}[i]$ to be the highest node $w$ with $\max(w) = i$). Let $j$ be the the leftmost occurrence of $\mathsf{pathlabel}(u)$ to the right of $i$. Then, the longest closed factor starting at position $i$ is $\mathsf{X}[i..j + |\mathsf{pathlabel}(u)| - 1]$ (this position $j$ is found by the KMP algorithm as in Lemma 1).

- If $u$ is the root, then it turns out that $i$ is the rightmost occurrence of character $\mathsf{X}[i]$ in $\mathsf{X}$. Hence, the longest closed factor starting at position $i$ is $\mathsf{X}[i]$.

The thing we have not shown is that $|\mathsf{pathlabel}(u)| = b_i$. This is indeed the case, since the set of occurrences of $\mathsf{X}[i..j + |\mathsf{pathlabel}(u)| - 1]$ (i.e., leaves in the subtree corresponding to the string) is equivalent to that of $\mathsf{pathlabel}(v)$, any substring starting at $i$ that is longer than $|\mathsf{pathlabel}(u)|$ does not occur to the right of $i$ and thus $b_i$ cannot be any longer. Hence $|\mathsf{pathlabel}(u)| = b_i$. (See also Figure 2).

Clearly $v = \mathsf{P}[i]$ can be retrieved in $\mathrm{O}(1)$ time for a given $i$, and then $u = \mathsf{parent}(v)$ can be obtained in $\mathrm{O}(1)$ time from $v$. This completes the proof. $\square$

The main result of this section follows is the following theorem.

**Theorem 1** *Given a string* $\mathsf{X}[1..n]$ *over an integer alphabet, the* longest closed factorization $\mathsf{LCF}(\mathsf{X}) = (\mathsf{G}_0, \mathsf{G}_1, \ldots, \mathsf{G}_k)$ *of* $\mathsf{X}$ *can be computed in* $\mathrm{O}(n)$ *time and space.*

PROOF. $\mathsf{G}_0 = \varepsilon$ by definition and so does not need to be computed. We compute the other $\mathsf{G}_j$ in ascending order of $j = 1, \ldots, k$. Let $s_i$ be the beginning position of $\mathsf{G}_i$ in $\mathsf{X}$, i.e., $s_1 = 1$ and $s_i = |\mathsf{G}_1 \cdots \mathsf{G}_{i-1}| + 1$ for $1 < i \leq k$. We compute $\mathsf{G}_1$ in $\mathrm{O}(|\mathsf{G}_1|)$ time and space from $b_{s_1}$ using Lemma 1 and Lemma 2. Assume we have computed the first $j - 1$ factors $\mathsf{G}_1, \ldots, \mathsf{G}_{j-1}$ for any $1 \leq j < k - 1$. We then

6

compute $\mathsf{G}_j$ in $O(|\mathsf{G}_j|)$ time and space from $b_{s_j}$, again using Lemmas 1 and 2. Since $\sum_{j=1}^{k} |\mathsf{G}_j| = n$, the proof completes. □

The following is an example of how the algorithm presented in this section computes $\mathsf{LCF}(\mathsf{X})$ for a given string $\mathsf{X}$.

**Example 7** *Consider the running example string* $\mathsf{X} = \texttt{ababaacbbbcbcc\$}$*, and see Figure 1, which shows the suffix tree of* $\mathsf{X}$.

1. *We begin with node* $\mathsf{P}[1]$ *representing* $\texttt{ababaacbbbcbcc\$}$*, whose parent represents* $\texttt{aba}$*. Hence we get* $b_1 = |\texttt{aba}| = 3$*. We run the KMP algorithm with pattern* $\texttt{aba}$ *and find the first factor* $\mathsf{G}_1 = \texttt{ababa}$.
2. *We then check node* $\mathsf{P}[6]$ *representing* $\texttt{a}$*. Since its parent is the root, we get* $b_2 = 0$ *and therefore the second factor is* $\mathsf{G}_2 = \texttt{a}$.
3. *We then check node* $\mathsf{P}[7]$ *representing* $\texttt{cbbbcbcc\$}$*, whose parent represents* $\texttt{cb}$*. Hence we get* $b_3 = |\texttt{cb}| = 2$*. We run the KMP algorithm with pattern* $\texttt{cb}$ *and find the third factor* $\mathsf{G}_3 = \texttt{cbbbcb}$.
4. *We then check node* $\mathsf{P}[13]$ *representing* $\texttt{cc\$}$*, whose parent represents* $\texttt{c}$*. Hence we get* $b_4 = |\texttt{c}| = 1$*. We run the KMP algorithm with pattern* $\texttt{c}$ *and find the fourth factor* $\mathsf{G}_4 = \texttt{cc}$.
5. *We finally check node* $\mathsf{P}[15]$ *representing* $\texttt{\$}$*. Since its parent is the root, we get* $b_5 = 0$ *and therefore the fifth factor is* $\mathsf{G}_5 = \texttt{\$}$.

*Consequently, we obtain* $\mathsf{LCF}(\mathsf{X}) = (\texttt{ababa}, \texttt{a}, \texttt{cbbbcb}, \texttt{cc}, \texttt{\$})$*, which coincides with Example 1.*


## 4. Longest Closed Factor Array

A natural extension of the problem in the previous section is to compute the longest closed factor starting at *every* position in $\mathsf{X}$ in linear time — not just those involved in the factorization. Formally, we would like to compute the longest closed factor array of $\mathsf{X}$, i.e., an array $\mathsf{LNG}[1..n]$ of integers such that $\mathsf{LNG}[i] = \ell$ if and only if $\ell$ is the length of the longest closed factor starting at position $i$ in $\mathsf{X}$.

Our algorithm for closed factorization computes the longest closed factor starting at a given position in time proportional to the factor's length, and so does not immediately provide a linear time algorithm for computing $\mathsf{LNG}$; indeed, applying the algorithm naïvely at each position would take $O(n^2)$ time to compute $\mathsf{LNG}$. In what follows, we present a more efficient solution.

**Theorem 2** *Given a string* $\mathsf{X}[1..n]$ *over an integer alphabet, the* longest closed factor array *of* $\mathsf{X}$ *can be computed in* $O(n \frac{\log n}{\log \log n})$ *time and* $O(n)$ *space.*

PROOF. We extend the data structure of the last section to allow $\mathsf{LNG}$ to be computed in $O(n \frac{\log n}{\log \log n})$ time and $O(n)$ space. The main change is to replace the KMP algorithm scanning in the first algorithm with a data structure that allows us to find the end of the closed factor in time independent of its length.

We first preprocess the suffix array $\mathsf{SA}$ for *range successor* queries, building the data structure of Yu, Hon and Wang [8]. A range successor query $\mathsf{rsq}_{\mathsf{SA}}(s, e, k)$ returns, given a range $[s, e] \subseteq [1, n]$, the smallest value $x \in \mathsf{SA}[s..e]$ such that $x > k$, or null if there is no value larger than $k$ in $\mathsf{SA}[s..e]$. Yu et al.'s data structure allows range successor queries to be answered in $O(\frac{\log n}{\log \log n})$ time each, takes $O(n)$ space, and $O(n \frac{\log n}{\log \log n})$ time to construct.

Now, to compute the longest closed factor starting at a given position $i$ in $\mathsf{X}$ (i.e. to compute $\mathsf{LNG}[i]$) we do the following. First we compute $b_i$, the length of the

border of the longest closed factor starting at $i$, in $O(1)$ time using Lemma 2. Recall that in the process of computing $b_i$ we determine the node $u$ having $\mathsf{pathlabel}(u) = \mathsf{X}[i..i + b_i - 1]$. To determine the end of the closed factor we must find the smallest $j > i$ such that $\mathsf{X}[j..j + b_i - 1] = \mathsf{X}[i..i + b_i - 1]$. Observe that $j$, if it exists, is precisely the answer to $\mathsf{rsq}_{\mathsf{SA}}(u.s, u.e, i)$. (See also the left diagram of Figure 2. Assuming that the leaves in the subtree rooted at $u$ are sorted in the lexicographical order, the leftmost and rightmost leaves in the subtree correspond to the $u.s$-th and $u.e$-th entries of $\mathsf{SA}$, respectively. Hence, $j = \mathsf{rsq}_{\mathsf{SA}}(u.s, u.e, i)$). For each $\mathsf{LNG}[i]$ we spend $O(\frac{\log n}{\log \log n})$ time and so overall the algorithm takes $O(n\frac{\log n}{\log \log n})$ time. The space requirement is clearly $O(n)$. $\qquad\square$

We note that recently Navarro and Neckrich [9] described range successor data structures with faster $O(\sqrt{\log n})$-time queries, but straightforward construction takes $O(n \log n)$ time [10], so overall this does not improve the runtime of our algorithm.

Given the longest closed factor array $\mathsf{LNG}$ of string $\mathsf{X}$, we can compute the maximal closed factor array $\mathsf{MAX}$ of $\mathsf{X}$ in linear time, by a simple scan on $\mathsf{LNG}$. Thus the next corollary follows from Theorem 2.

**Corollary 1** *Given a string* $\mathsf{X}[1..n]$ *over an integer alphabet, the* maximal closed factor array *of* $\mathsf{X}$ *can be computed in* $O(n\frac{\log n}{\log \log n})$ *time and* $O(n)$ *space.*

## 5. Shortest Closed Factors

In this section, we consider the problems related to shortest closed factors of length at least 2 in a given string. The following simple lemma is a key to our algorithms.

**Lemma 3** *Any shortest closed factor of length at least 2 of string* $\mathsf{X}$ *has a unique border of length 1.*

PROOF. For any position $i$ in $\mathsf{X}$, let $\mathsf{X}[i..j]$ be the shortest closed factor of length at least 2 starting at position $i$. Since $j - i + 1 \geq 2$ and $\mathsf{X}[i..j]$ is the shortest closed factor starting at position $i$, $\mathsf{X}[i..j]$ must have a unique border of length at least 1. Assume, for the sake of contradiction, that the border of $\mathsf{X}[i..j]$ is of length $k \geq 2$. If $\mathsf{X}[j - k + 1..j - 1]$ occurs exactly twice in $\mathsf{X}[i..j - 1]$, then $\mathsf{X}[i..j - 1]$ is also closed. On the other hand, if it occurs twice or more, let $\mathsf{X}[p..p + k - 2]$, with $i < p < j - k + 1$, be the first occurrence in $\mathsf{X}[i + 1..j - 2]$. Then $\mathsf{X}[i..p + k - 2]$ is also closed. However, both cases contradict that $\mathsf{X}[i..j]$ is the shortest closed factor starting at position $i$. Hence the unique border of $\mathsf{X}[i..j]$ is of length 1, namely, a single character. $\qquad\square$

We now have the following theorem.

**Theorem 3** *Given a string* $\mathsf{X}[1..n]$ *over an integer alphabet of size* $\sigma$, *the* shortest closed factor array $\mathsf{SHT}$ *can be computed in* $O(n)$ *time and* $O(\sigma)$ *space,*

PROOF. For any position $i$ in $\mathsf{X}$, let $\mathsf{X}[i] = c$. Let $j$ be the smallest position in $\mathsf{X}$ such that $i < j$ and $\mathsf{X}[j] = c$, if it exists. It follows from Lemma 3 that $\mathsf{SHT}[i] = j - i + 1$ if $j$ exists, and $\mathsf{SHT}[i] = null$ otherwise. Hence $\mathsf{SHT}$ can be computed in $O(n)$ time by a simple scan on $\mathsf{X}$, using $O(\sigma)$ space. $\qquad\square$

The next corollary is immediate from Theorem 3.

**Corollary 2** *Given a string* $\mathsf{X}[1..n]$ *over an integer alphabet of size* $\sigma$, *the* shortest closed factorization $\mathsf{SCF}(\mathsf{X}) = (\mathsf{F}_0, \mathsf{F}_1, \ldots, \mathsf{F}_h)$ *can be computed in* $\mathrm{O}(n)$ *time and* $\mathrm{O}(\sigma)$ *space if it exists.*

Now we consider the problem of computing, for each position $i$ of a given string, the length of a shortest closed factor containing $i$. The following lemma is useful to solve this problem.

**Lemma 4 (Min-Variant of the Manhattan Skyline Problem [11])** *Let $S$ be a set of* $\mathrm{O}(n)$ *integer subintervals of an interval* $[1..n]$ *each with an associated non-negative integer* height *of size* $\mathrm{O}(n)$; *the height value for interval* $[i..j] \in S$ *is denoted* $height([i..j])$. *Then, we can compute a table $T$ such that, for any $1 \leq t \leq n$,* $T[t] = \min\{height([i..j]) \mid t \in [i..j], [i..j] \in S\}$ *in* $\mathrm{O}(n)$ *time and space.*

**Theorem 4** *Given a string* $\mathsf{X}[1..n]$ *over an integer alphabet, the* minimal closed factor array $\mathsf{MIN}$ *can be computed in* $\mathrm{O}(n)$ *time and space.*

PROOF. We first compute the shortest closed factor array $\mathsf{SHT}$ of $\mathsf{X}$ in $\mathrm{O}(n)$ time and $\mathrm{O}(\sigma)$ space by Theorem 3. For each $1 \leq i \leq n$, we associate $\mathsf{SHT}[i]$ to the interval $[i..i + \mathsf{SHT}[i] - 1]$ as its height. Now, the problem of computing the $\mathsf{MIN}$ array is equivalent to the min-variant of the Manhattan skyline problem (defined above). Hence, the $\mathsf{MIN}$ array can be computed in $\mathrm{O}(n)$ time and space by Lemma 4. $\qquad\square$

## 6. Concluding Remarks

We have considered several problems on closed factors here, and many others remain. For example, how efficiently can one compute all the closed factors in a string (or, say, the closed factors that occur at least $k$ times)? Relatedly, what is the maximum number of closed factors a string can contain? What is the average number?

One also wonders if the longest closed factor array can be computed in linear time, by somehow avoiding range successor queries. Another open question is whether the maximal closed factor array can be computed in linear time, without using the longest closed factor array.

## References

[1] G. Fici, A classification of trapezoidal words, in: Proc. 8th International Conference Words 2011 (WORDS 2011), Electronic Proceedings in Theoretical Computer Science 63, 2011, pp. 129–137, see also http://arxiv.org/abs/1108.3629v1.

[2] G. Badkobeh, G. Fici, Z. Lipták, A note on words with the smallest number of closed factors, http://arxiv.org/abs/1305.6395 (2013).

[3] G. Badkobeh, H. Bannai, K. Goto, T. I, C. S. Iliopoulos, S. Inenaga, S. J. Puglisi, S. Sugimoto, Closed factorization, in: Proc. PSC 2014, 2014, pp. 162–168.

[4] U. Manber, G. W. Myers, Suffix arrays: a new method for on-line string searches, SIAM Journal on Computing 22 (5) (1993) 935–948.

[5] P. Weiner, Linear pattern matching, in: IEEE 14th Annual Symposium on Switching and Automata Theory, IEEE, 1973, pp. 1–11.

[6] M. Farach, Optimal suffix tree construction with large alphabets, in: Proceedings of the 38th Annual Symposium on Foundations of Computer Science, 1997, pp. 137–143.

[7] D. E. Knuth, J. H. Morris, V. R. Pratt, Fast pattern matching in strings, SIAM J. Comput. 6 (2) (1977) 323–350.

[8] C.-C. Yu, W.-K. Hon, B.-F. Wang, Improved data structures for the orthogonal range successor problem, Computational Geometry 44 (3) (2011) 148–159.

[9] G. Navarro, Y. Nekrich, Sorted range reporting, in: Proc. Scandinavian Workshop on Algorithm Theory, LNCS 7357, Springer, 2012, pp. 271–282.

[10] G. Navarro, Y. Nekrich, personal Communication.

[11] M. Crochemore, C. S. Iliopoulos, M. Kubica, J. Radoszewski, W. Rytter, T. Walen, Extracting powers and periods in a word from its runs structure, Theoretical Computer Science 521 (2014) 29–41.