

Computing Longest Common Substrings Via Suffix Arrays

Maxim A. Babenko and Tatiana A. Starikovskaya

Moscow State University
max@adde.math.msu.su,
tat.starikovskaya@gmail.com

Abstract. Given a set of N strings $A = \{\alpha_1, \dots, \alpha_N\}$ of total length n over alphabet Σ one may ask to find, for each $2 \leq K \leq N$, the longest substring β that appears in at least K strings in A . It is known that this problem can be solved in $O(n)$ time with the help of suffix trees. However, the resulting algorithm is rather complicated (in particular, it involves answering certain least common ancestor queries in $O(1)$ time). Also, its running time and memory consumption may depend on $|\Sigma|$.

This paper presents an alternative, remarkably simple approach to the above problem, which relies on the notion of suffix arrays. Once the suffix array of some auxiliary $O(n)$ -length string is computed, one needs a simple $O(n)$ -time postprocessing to find the requested longest substring. Since a number of efficient and simple linear-time algorithms for constructing suffix arrays has been recently developed (with constant not depending on $|\Sigma|$), our approach seems to be quite practical.

1 Introduction

Consider the following problem:

(LCS) *Given a collection of N strings $A = \{\alpha_1, \dots, \alpha_N\}$ over alphabet Σ find, for each $2 \leq K \leq N$, the longest string β that is a substring of at least K strings in A .*

It is known as a generalized version of the *Longest Common Substring* (LCS) problem and has a plenty of practical applications, see [Gus97] for a survey.

Even in the simplest case of $N = K = 2$ a linear-time algorithm is not easy. A standard approach is to construct the so-called *generalized suffix tree* T (see [Gus97]) for $\alpha_1\$1$ and $\alpha_2\$2$, which is a compacted symbol trie that captures all the substrings of $\alpha_1\$1$, $\alpha_2\$2$. Here $\$i$ are special symbols (called *sentinels*) that are distinct and do not appear in α_1 and α_2 . Then, nodes of T are examined in a bottom-up fashion and those having sentinels of both types in their subtrees are listed. Among these nodes of T let us choose a node v with the largest *string depth* (which is the length of the string obtained by reading letters along the path from root to v). The string that corresponds to v in T is the answer. See [Gus97] for more details.

In practice, the above approach is not very efficient since it involves computing T . Several linear-time algorithms for the latter task are known (possibly, the most famous one is due to Ukkonen [Ukk95]). However, suffix trees are still not very convenient. They do have linear space bound but the hidden constants can be pretty large. Most of modern algorithms for computing suffix trees have the time bound of $O(n \log |\Sigma|)$ (where n denotes the length of a string). Hence, their running time depends on $|\Sigma|$. Moreover, achieving this time bound requires using balanced search trees to store arcs. The latter data structures further increase constants in both time- and space-bounds making these algorithms rather impractical. Other options include employing hashables or assuming that $|\Sigma|$ is small and using direct addressing to access arcs leaving each node. These approaches have their obvious disadvantages.

If one assumes that N and K are arbitrary then additional complications arise. Now we are interested in finding the deepest (in sense of string depth) node v in T such that the tree rooted at v contains sentinels of at least K distinct kinds. Moreover, this routine should run in parallel for all possible values of K and take linear time. This seems to be an involved task. A possible solution requires answering *Least Common Ancestor* (LCA) queries on T in $O(1)$ time, e.g. using a method from [BFC00]. Reader may refer to [Gus97] for a complete outline.

In this paper we present an alternative approach that is based on the notion of *suffix arrays*. The latter were introduced by Manber and Myers [MM90] in an attempt to overcome the issues that are inherent to suffix trees. The *suffix array* (SA) of string α having length n is merely an array of n integers that indicate the lexicographic order of non-empty suffixes of α (see Section 2 for a precise definition). Its simplicity and compactness make it an extremely useful tool in modern text processing. Originally, an $O(n \log n)$ -time algorithm for constructing SA was suggested [MM90]. This algorithm is not very practical. Subsequently, much simpler and faster algorithms for computing SA were developed. We particularly mention an elegant approach of Kärkkäinen and Sanders [KS03]. A comprehensive practical evaluation of different algorithms for constructing SA is given in [PST07].

We present two algorithms. The first one, which is simpler, assumes that parameter K is fixed. It first builds an auxiliary string α by concatenating strings α_i and intermixing them with sentinels $\$i$ ($1 \leq i \leq N$) and then constructs the suffix array for string α . Also, an additional *LCP array* is constructed. Finally, a sliding window technique is applied to these arrays to obtain the answer. Altogether, the running time is linear and does not depend on $|\Sigma|$.

The second algorithm deals with all possible values of K simultaneously. Its initial stage is similar: string α , suffix array for α , and LCP array are constructed. Then, *Cartesian tree* (CT) is constructed from LCP array. Finally, a certain postprocessing aimed to count the number of distinct types of nodes appearing in subtrees of CT is used. It should be noticed that this postprocessing does not require answering any least common ancestor queries.

The paper is organized as follows. Section 2 gives a formal background, introduces useful notation and definitions. It also explains the notion of suffix arrays

and indicates how an auxiliary LCP array is constructed in linear time. Section 3 presents the algorithm for the case of fixed K . Finally, Section 4 considers the general case when the value of K is not fixed.

2 Preliminaries

We shall start with a number of definitions first. In what follows we assume that a finite non-empty set Σ (called an *alphabet*) is fixed. The elements of Σ are *letters* or *symbols*. A finite ordered sequence of letters (possibly empty) is called a *string*.

We assume the usual RAM model of computation [AUH74]. Letters are treated just as integers in range $\{1, \dots, |\Sigma|\}$, so one can compare any pair of them in $O(1)$ time. This *lexicographic* order on Σ is linear and can be extended in a standard way to the set of strings in Σ . We write $\alpha < \beta$ to denote that α lexicographically precedes β ; similarly for other relation signs.

We usually use Greek symbols to denote strings. Letters in a string are numbered starting from 1, that is, for a string α of length k its letters are $\alpha[1], \dots, \alpha[k]$. The length k of α is denoted by $|\alpha|$. The *substring* of α from position i to position j (inclusively) is denoted by $\alpha[i..j]$. Also, if $i = 1$ or $j = |\alpha|$ then these indices are omitted from the notation and we write just $\alpha[..j]$ and $\alpha[i..]$. String $\beta = \alpha[..j]$ is called a *prefix* of α . Similarly, if $\beta = \alpha[i..]$ then β is called a *suffix* of α . For a set of strings S let $\text{lcp}(S)$ denote the longest common prefix of all strings in S .

Recall that A stands for the collection of the input strings α_i . We start with an almost trivial observation:

Proposition 1. *Let $B = \{\beta_1, \dots, \beta_m\} \subseteq A$ be an arbitrary subset of A obeying $m \geq K$. Let γ_i be an arbitrary suffix of β_i for each $1 \leq i \leq m$. Solving (LCS) amounts to computing the longest string among*

$$\text{lcp}(\gamma_1, \dots, \gamma_m)$$

where maximum is taken over all possible choices of subsets B and suffixes $\{\gamma_i\}$.

Let us combine the strings in A as follows:

$$\alpha = \alpha_1 \$1 \alpha_2 \$2 \dots \alpha_N \$N \tag{1}$$

Here $\$_i$ are pairwise distinct *sentinel* symbols not appearing in strings of A . We assume that these sentinels are lexicographically smaller than other (normal) symbols. The lexicographic order between sentinels is not important.

String α captures all needed information about set A . For each index i ($1 \leq i \leq N$) and a position j in α_i ($1 \leq j \leq |\alpha_i|$) one may consider the *corresponding* position $p(i, j)$ in α :

$$p(i, j) := \sum_{k=1}^{i-1} (|\alpha_k| + 1) + j$$

	suffixes	SA	sorted suffixes	lcp
1	mississippi	11	i	1
2	ississippi	8	ippi	1
3	ssissippi	5	issippi	4
4	sissippi	2	ississippi	0
5	issippi	1	mississippi	0
6	ssippi	10	pi	1
7	sippi	9	ppi	0
8	ippi	7	sippi	2
9	ppi	4	sissippi	1
10	pi	6	ssippi	3
11	i	3	ssissippi	

Fig. 1. String `mississippi`, its suffixes, and the corresponding suffix and LCP arrays

Positions in α of the form $p(i, j)$ are called *essential*; the remaining positions (those containing sentinels) are called *unessential*.

Let us employ the following metaphor: for each $1 \leq i \leq N$ and $1 \leq j \leq |\alpha_i|$ we say that position $p(i, j)$ is of *type* i (it corresponds to the i -th string). Remaining (unessential) positions k in α are said to be of type 0.

Now taking into account the properties of sentinels one can easily derive the following claim from Proposition 1:

Proposition 2. *Let $P = \{p_1, \dots, p_m\}$ be an arbitrary set of essential positions in α such that elements of P are of at least K distinct types. Solving (LCS) amounts to computing the longest string among*

$$lcp(\alpha[p_1..], \dots, \alpha[p_m..])$$

where maximum is taken over all possible choices of P .

This does not seem very promising at the first glance. However, the longest common prefix computation exhibits a nice structure when it is applied to suffixes of a fixed string (in our case, string α).

To explain this structure we first introduce the notion of suffix arrays. Let ω be an arbitrary string of length n . Consider its non-empty suffixes

$$\omega[1..], \omega[2..], \dots, \omega[n..]$$

and order them lexicographically. Let $SA(i)$ denote the starting position of the suffix appearing on the i -th place ($1 \leq i \leq n$) in this order:

$$\omega[SA(1)..] < \omega[SA(2)..] < \dots < \omega[SA(n)..]$$

Clearly, SA is determined uniquely since all suffixes of ω are distinct. An example is depicted in Fig. 1.

Since SA is a permutation of $\{1, \dots, n\}$ there must be an inverse correspondence. We denote it by *rank*; that is, *rank* is also a permutation of $\{1, \dots, n\}$ and

$$SA(rank(i)) = i \quad \text{holds for all } 1 \leq i \leq n.$$

Historically, the first algorithm to compute SA was due to Manber and Myers [MM90]; this algorithm takes $O(n \log n)$ time. Currently, simple linear-time algorithms for this task are known, see [PST07] for a list. The latter linear time bounds do not depend on $|\Sigma|$.

However, knowing SA is not enough for our purposes. We also have to precompute the lengths of longest common prefixes for each pair of consequent suffixes (with respect to the order given by SA). More formally,

$$lcp(i) := |lcp(\omega[SA(i)..], \omega[SA(i+1)..])| \quad \text{for all } 1 \leq i < n.$$

This gives rise to array lcp of length $n - 1$; we call it the *LCP array* of ω . The latter array not only enables to answer LCP queries for consequent (w.r.t. SA) suffixes of ω but also carries enough information to answer *any* such query. Formally [MM90]:

Lemma 1. *For each pair $1 \leq i < j \leq n$ one has*

$$|lcp(\omega[SA(i)..], \omega[SA(j)..])| = \min_{i \leq k < j} lcp(k)$$

Knowing the suffix array, LCP array may be constructed in $O(n^2)$ time by a brute-force method. However, an elegant modification ([KLA⁺01], see also [CR03]) allows to compute the longest common prefix for a pair of consequent suffixes in $O(1)$ amortized time. The key is to compute these values in a particular order, namely

$$lcp(rank(1)), lcp(rank(2)), \dots, lcp(rank(n))$$

The efficiency of this approach relies on the following fact [KLA⁺01]:

Lemma 2. *$lcp(rank(i+1)) \geq lcp(rank(i)) - 1$ for each $1 \leq i < n$ such that $rank(i) < n$ and $rank(i+1) < n$.*

Hence, when computing the value of $lcp(rank(i+1))$ one can safely skip $lcp(rank(i)) - 1$ initial letters of $\omega[i+1..]$ and $\omega[SA(rank(i+1)+1)..]$. This easily implies the required linear time bound for the whole computation.

3 Fixed K : Sliding Window

We now proceed by describing the algorithm that solves (LCS) for a fixed value of K . There are several reasons for considering this case separately. Firstly, this problem also seems natural and the resulting approach is somewhat simpler. Secondly, the techniques that we develop for this special case turn out to be useful for the general problem.

The algorithm works as follows. It first combines the input strings into string α of length L (see (1)) and invokes the suffix array computation algorithm thus obtaining the suffix array SA for α . It also constructs array lcp in $O(L)$ time as described in Section 2.

Refer to Proposition 2 and consider an arbitrary set of essential positions $P = \{p_1, \dots, p_m\}$ that contains positions of at least K distinct types. Replace these positions with ranks by putting $r_i := \text{rank}(p_i)$ and, thus, forming the set $R = \{r_1, \dots, r_m\}$. Lemma 1 implies the equality

$$|\text{lcp}(\alpha[p_{1..}], \dots, \alpha[p_{m..}])| = \min_{R^- \leq j < R^+} \text{lcp}(j)$$

where $R^- := \min R$ and $R^+ := \max R$.

Let us consider a segment $\Delta \subseteq [1, L]$ and call it K -good if for $i \in \Delta$ positions $SA(i)$ are of at least K distinct essential types. This enables us to restate Proposition 2 as follows:

Proposition 3. *The length of the longest common substring that appears in at least K input strings is equal to*

$$\max_{\Delta} \min_{\Delta^- \leq j < \Delta^+} \text{lcp}(j)$$

where $\Delta = [\Delta^-, \Delta^+]$ ranges over all K -good segments.

This formula is already an improvement (compared to Proposition 2) since it only requires to consider a polynomial number of possibilities. Moreover, we shall indicate how maximum in Proposition 3 can be found in $O(L)$ time.

To this aim, note that if a K -good segment Δ is already considered then any $\Delta' \supset \Delta$ cannot give us a bigger value of minimum. For each i one may consider the segment $\Delta_i = [\Delta_i^-, \Delta_i^+]$ obeying the following properties:

- Δ_i starts at position i ;
- Δ_i is K -good;
- Δ_i is the shortest segment obeying the above conditions.

Here index i ranges over $[1, L_0]$, where $L_0 \leq L$ is the smallest integer such that segment $[L_0 + 1, L]$ is not K -good.

Note that due to our assumption that sentinels are strictly less than normal letters, the first N elements of SA correspond to unessential positions occupied by the sentinels. The algorithm does not need to consider these positions and only examines segments $\Delta_{N+1}, \dots, \Delta_{L_0}$.

Put

$$w(i) := \min_{\Delta_i^- \leq j < \Delta_i^+} \text{lcp}(j) \quad \text{for all } N < i \leq L_0.$$

and consider the sequence:

$$w(N + 1), w(N + 2), \dots, w(L_0) \tag{2}$$

Once the maximum among (2) is found, the problem is solved. We construct a pipeline with the first stage computing the sequence of segments

$$\Delta_{N+1}, \Delta_{N+2}, \dots, \Delta_{L_0}$$

and the second stage calculating the respective minima (2)

Put $\Delta_i = [\Delta_i^-, \Delta_i^+]$. The first stage works as follows. It initially sets $\Delta_{N+1}^- := N + 1$ and then advances the right endpoint Δ_{N+1}^+ until getting a K -good segment. Then, on each subsequent iteration i it puts $\Delta_i^- := i$, $\Delta_i^+ := \Delta_{i-1}^+$ and again advances the right endpoint Δ_i^+ until Δ_i becomes K -good. (In case, no such segment can be extracted, it follows that $i > L_0$, so the end is reached.)

Lemma 3. Δ_i is the shortest K -good segment starting at i .

Proof. We claim that for any K -good segment $[i, j]$ one has $j \geq \Delta_{i-1}^+$. Indeed, suppose towards contradiction that $j < \Delta_{i-1}^+$. Since $[i, j]$ is K -good then so is $[i - 1, j]$. The latter, however, contradicts the minimality of Δ_{i-1} .

To test in $O(1)$ time if the current candidate forms a K -good segment the algorithm maintains an array of counters $c(1), \dots, c(N)$. For each $N < j \leq L$ put $t(j)$ to be the type of position $SA(j)$ in α (recall that all these positions are essential). For each index i ($1 \leq i \leq N$) the entry $c(i)$ stores the number of positions j in the current segment such that $t(j) = i$. Also, the number of non-zero entries of c (denoted by $npos$) is maintained.

Initializing c and $npos$ for Δ_{N+1} is trivial. Then, when the algorithm puts $\Delta_i^- = \Delta_{i-1}^- + 1$ it decrements the entry of c that corresponds to position $i - 1$ (which has just been removed from the window) and adjusts $npos$, if necessary. Similarly, when the current segment is extended to the right, certain entries of c are increased and $npos$ is adjusted. To see whether the current segment is K -good one checks if $npos \geq K$. This completes the description of the first stage of the pipeline. Note that it totally takes $O(L)$ time.

The second stage aims to maintain the values (2) dynamically. This is done by the following (possibly folklore) trick. Consider a queue Q that, at any given moment i , holds the sequence of keys

$$lcp(\Delta_i^-), lcp(\Delta_i^- + 1), \dots, lcp(\Delta_i^+ - 1)$$

Increasing index i the algorithm dequeues value $lcp(\Delta_i^-)$ from the head of Q (to account for the increase of Δ_i^-) and then enqueues some (possibly none) additional values to the tail of Q (to account for the increase of Δ_i^+ , if any). The total number of these queue operations is $O(L)$.

We describe a method for maintaining the minimum of keys in Q under insertions and removals and serving each such request in $O(1)$ amortized time. A queue Q that holds a sequence (q_1, \dots, q_m) may be simulated by a pair of stacks S^1 and S^2 . A generic configuration of these stacks during this simulation is as follows (here $0 \leq s \leq m$):

$$\begin{aligned} S^1 &= (q_s, q_{s-1}, \dots, q_1) \\ S^2 &= (q_{s+1}, q_{s+2}, \dots, q_m) \end{aligned} \tag{3}$$

Here stack elements are listed from bottom to top. Initially Q is empty, hence so are S^1 and S^2 . To enqueue a new key x (which becomes q_{m+1}) to Q one pushes x onto S^2 . This takes $O(1)$ time. To dequeue q_1 from Q consider two cases. If $s > 0$

then S^1 is non-empty; pop the top element q_1 from S^1 . Otherwise, one needs to transfer elements from S^2 to S^1 . This is done by popping elements from S^2 one after another and simultaneously pushing them onto S^1 (in the same order). By (3) these operations preserve the order of keys in Q . Once they are complete, S^1 becomes non-empty and the first case applies. To estimate the running time note that any enqueued element may participate in an S^2 -to- S^1 transfer at most once. Hence, an amortized bound of $O(1)$ follows.

Our algorithm simulates Q via S^1 and S^2 , as explained above. Each S^i is additionally augmented to maintain the minimum among the keys it contains. This is achieved by keeping minima m^1 , m^2 and a pair of auxiliary stacks M^1 , M^2 . When a new key x is pushed to S^i the algorithm saves the previous minimum m^i in M^i and updates m^i by $m^i := \min(m^i, x)$. When an element is popped from S^i the algorithm also pops m^i from M^i .

Altogether these manipulations with Q , S^i , M^i , and m^i take time that is proportional to the number of operations applied to Q . The latter is known to be $O(L)$. Hence, the algorithm computes the sequence of minima (2) and chooses the maximum (call it $M(K)$) among these values in $O(L)$ time, as claimed.

Let the above maximum be attained by a segment $\Delta = [\Delta^-, \Delta^+] \subseteq [N+1, L]$. Suppose that position $SA(\Delta^-)$ in string α corresponds to some position j in some input string α_i . Now the desired longest common substring is $\alpha_i[j..j+M(K)-1]$.

4 Arbitrary K : Cartesian Tree

Now we go back to the original problem (LCS) and no longer assume that the value of K is given in advance. Similarly to Section 3, we start by constructing string α , suffix array SA of α , and the corresponding LCP array.

The cornerstone of the algorithm is a novel postprocessing, which combines the above information. Firstly, we need to overcome the following technical issue. Recall from the previous section that $t(j) \in \{1, \dots, N\}$ denotes the type of position $SA(j)$ in string α ($N < j \leq L$). Also, a segment Δ is called K -good if $t(j)$ gives at least K distinct values while j ranges over Δ . According to Proposition 3 to estimate the length of the common substring corresponding to segment Δ one needs to compute $\min_j lcp(j)$ where j ranges over Δ *without its right endpoint*.

To simplify the matters we construct a new pair of arrays lcp' and t' (whose elements are numbered starting from 1) by intermixing elements of lcp and t with artificial values as follows:

$$\begin{aligned} lcp' &:= (\infty, lcp(N+1), \infty, lcp(N+2), \dots, lcp(L-1), \infty) \\ t' &:= (t(N+1), 0, t(N+2), 0, \dots, 0, t(L)) \end{aligned}$$

Now for arrays lcp' and t' it is clear that the same segment Δ should be used both for calculating the number of distinct non-zero values among $t'(j)$, $j \in \Delta$ and the respective minima $\min(lcp'(j) : j \in \Delta)$.

We remind the reader the notion of *Cartesian trees* (see, e.g., [BFC00]). Let $A = (a_1, \dots, a_n)$ be an arbitrary (possibly empty) sequence of items, and let

each a_i be assigned a real number $key(a_i)$. We associate a tree $CT(A)$ with A by the following rules:

- if $n = 0$ then $CT(A)$ is the *null tree* (having no nodes);
- if $n > 0$ then let a_i denote an item in A with the smallest key $key(a_i)$ (the choice of i may not be unique); now $CT(A)$ is constructed by taking a_i as its root, recursively constructing trees $CT_L := CT(a_1, \dots, a_{i-1})$ and $CT_R := CT(a_{i+1}, \dots, a_n)$, and attaching CT_L to a_i as its left child and CT_R as its right child.

Proposition 4. *Consider a pair $1 \leq i \leq j \leq n$ and let $lca(a_i, a_j)$ denote the least common ancestor of nodes a_i and a_j in $CT(A)$. Then*

$$\min_{i \leq k \leq j} key(a_k) = key(lca(a_i, a_j)).$$

Let L' denote the length of lcp' and t' . We construct a Cartesian tree CT taking positions $i \in \{1, \dots, L'\}$ as nodes and using values $lcp'(i)$ as keys. For a node v in CT let CT_v denote the subtree rooted at v and $z(v)$ — the number of distinct non-zero values of $t'(u)$ for $u \in CT_v$.

The next reformulation of (LCS) is an immediate consequence of Proposition 4:

Proposition 5. *The length of the longest common substring that appears in at least K input strings is equal to $\max_v lcp'(v)$, where v ranges over all nodes of CT such that $z(v) \geq K$.*

This provides us with the following two problems: how to construct CT from lcp' values and how to compute z values. The first task is easily solvable in $O(L)$ time (see, e.g., [BFC00]); we shall describe this method below.

The second task is more involved. Instead of working directly with values of z consider an integer-valued function ζ on the nodes of CT such that

$$z(v) = \sum_{u \in CT_v} \zeta(u).$$

Clearly, ζ is uniquely determined by z and vice versa. We construct ζ and CT incrementally by scanning arrays lcp' and t' and adding, at the i -th step ($1 \leq i \leq L'$), node i with key $lcp'(i)$ and type $t'(i)$.

The algorithm maintains a stack $S = (S(1), \dots, S(m))$ that holds nodes of CT contained on the *rightmost path* (that is, the path formed by walking from the root of CT and taking the right child on each step until reaching a null reference). Here m is the length of the rightmost path, $S(1)$ is the root of CT (the bottom of the stack), and $S(m)$ is the last node on the rightmost path (the top of the stack).

Adding the first node 1 to CT is straightforward. Suppose that nodes $1, \dots, i$ are already added to CT and the corresponding values of ζ are computed. To insert node $i + 1$ into CT the algorithm iteratively pops a sequence of nodes $S(m), S(m - 1), \dots, S(l)$ from the stack as long as the last popped node $S(l)$

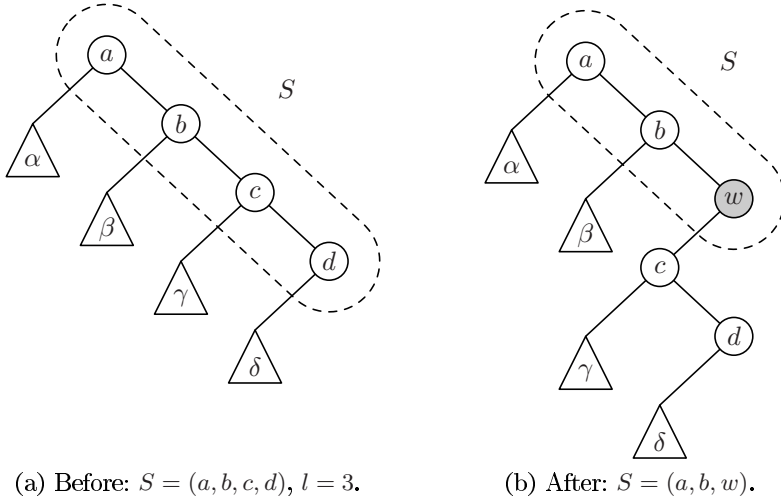


Fig. 2. Inserting a new node w into a Cartesian tree. Here $key(a) \leq key(b) < key(w) \leq key(c) \leq key(d)$.

obeys $lcp'(S(l)) \geq lcp'(i + 1)$. In particular, if $lcp'(S(m)) < lcp'(i + 1)$ then $l := m + 1$ and no nodes are popped. Tree CT is adjusted as follows (see Fig. 2 for an example):

- the right child of $i + 1$ is set to null;
- the left child of $i + 1$ is set to $S(l)$ if $l \leq m$ or to null if $l > m$;
- node $i + 1$ is added to the end of the rightmost path.

The total time that is necessary to perform these operations for node $i + 1$ is $O(2 + \Delta m)$, where Δm denotes the decrease of m . By amortization, this easily yields a linear time bound for the total process.

Now we need to explain how the values of ζ are changed by the insertion. If $t'(i + 1) = 0$ then it suffices to put $\zeta(i + 1) := 0$. Otherwise, let $t'(i + 1) \neq 0$. Clearly, we may only need to change values $\zeta(S(1)), \dots, \zeta(S(l - 1))$ (if any) and to initialize $\zeta(i + 1)$; other values of ζ correspond to subtrees of CT that are not affected by the insertion.

At any given moment we assign *ranks* to the nodes of CT by the following rule: for each $1 \leq i \leq m$ node $S(i)$ and all nodes in its left subtree are of rank i . For each possible type $j \in \{1, \dots, N\}$ let $max\text{-rank}(j)$ denote the maximum rank of a node in CT having type j . In case no node of type j exists, we put $max\text{-rank}(j) := -\infty$.

Two cases are possible. Firstly, suppose $l' := max\text{-rank}(t'(i + 1)) \geq l$. This means that prior to insertion of node $i + 1$ the first node of type $t'(i + 1)$ (with regard to the in-order traversal) was occurring in $CT_{S(l)}$. Then, the step is completed by putting $\zeta(i + 1) := 0$.

Secondly, let $l' < l$ (in particular, this case applies when $l' = -\infty$). The algorithm assigns $\zeta(i+1) := 1$ and decreases the value $\zeta(S(l'))$ by 1 (if $l' > -\infty$).

It is easy to see that these changes are correct and produce the required values of ζ . However, it remains to explain how the *max-rank* values are computed. To this aim, we need a data structure for maintaining an array *max-rank* = (*max-rank*(1), ..., *max-rank*(N)) of these values. In addition to the standard read and write requests, this array should also be capable of performing *trimming* as follows: given a value l put

$$\text{max-rank}(j) := \min(\text{max-rank}(j), l) \quad \text{for all } 1 \leq j \leq N.$$

This operation is invoked to adjust the maximum ranks each time node $S(l)$ is turned into a left child of node $i + 1$.

A possible implementation is based on keeping, for each index $1 \leq j \leq N$, an integer *timestamp last-write*(j) that keeps the moment of time when this entry was last updated. Also, instead of performing it directly, the algorithm maintains the timestamp *last-trim* of the latest trimming operation (together with the corresponding trimming parameter l). Now to get the actual value *max-rank*(j) one compares *last-write*(j) with *last-trim* to see if trimming applies to the currently stored value. With this implementation, each read access, write access or trimming takes $O(1)$ time.

Therefore, we can construct CT and compute the values of ζ in linear time. Then, the values of z are computed from ζ in a bottom-up fashion. Simultaneously, for each possible value of $z(i)$ we accumulate the largest value of $lcp'(i)$ and construct the array

$$\text{max-lcp}'(k) := \max(lcp'(i) : z(i) = k) \quad \text{for all } 2 \leq k \leq N.$$

The length of the longest common substring corresponding to a certain value of K is

$$M(K) := \max(\text{max-lcp}'(k) : k \geq K) \quad \text{for all } 2 \leq K \leq N. \quad (4)$$

Hence, all values $M(2), \dots, M(N)$ may be computed by an obvious recurrence.

The longest substrings themselves may also be easily extracted. For each node v of CT enumerate the nodes in CT_v via an in-order traversal and put u^- (resp. v^+) to be the first (resp. the last) node in CT_v such that $t'(v^-) \neq 0$ (resp. $t'(v^+) \neq 0$). Clearly, computing all nodes v^+ and v^- takes linear time.

Now consider a fixed value of K . Let the maximum in (4) be attained by some $k \geq K$. Next, let v denote a node in CT such that $z(v) = k$. Nodes v^- and v^+ give rise to a k -good segment $\Delta = [\Delta^-, \Delta^+] \subseteq [N + 1, L]$ obeying the equality $\min(lcp(j) : \Delta^- \leq j < \Delta^+) = M(K)$. The corresponding longest common substring is constructed from Δ in the same way as in Section 3.

Acknowledgements

The authors are thankful to the students of Department of Mathematical Logic and Theory of Algorithms (Faculty of Mechanics and Mathematics, Moscow

State University), and also to Maxim Ushakov and Victor Khimenko (Google Moscow) for discussions.

References

- [AUH74] Aho, A.V., Ullman, J.D., Hopcroft, J.E.: The design and analysis of computer algorithms. Addison-Wesley, Reading, MA (1974)
- [BFC00] Bender, M., Farach-Colton, M.: The LCA problem revisited, pp. 88–94 (2000)
- [CR03] Crochemore, M., Rytter, W.: Jewels of stringology. World Scientific Publishing Co. Inc., River Edge, NJ (2003)
- [Gus97] Gusfield, D.: Algorithms on strings, trees, and sequences: Computer science and computational biology. Cambridge University Press, New York, NY, USA (1997)
- [KLA⁺01] Kasai, T., Lee, G., Arimura, H., Arikawa, S., Park, K.: Linear-time longest-common-prefix computation in suffix arrays and its applications. In: Amir, A., Landau, G.M. (eds.) CPM 2001. LNCS, vol. 2089, pp. 181–192. Springer, Heidelberg (2001)
- [KS03] Kärkkäinen, J., Sanders, P.: Simple linear work suffix array construction. In: Baeten, J.C.M., Lenstra, J.K., Parrow, J., Woeginger, G.J. (eds.) ICALP 2003. LNCS, vol. 2719. Springer, Heidelberg (2003)
- [MM90] Manber, U., Myers, G.: Suffix arrays: a new method for on-line string searches. In: SODA 1990: Proceedings of the first annual ACM-SIAM symposium on Discrete algorithms. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, pp. 319–327 (1990)
- [PST07] Puglisi, S.J., Smyth, W.F., Turpin, A.H.: A taxonomy of suffix array construction algorithms. ACM Comput. Surv. 39(2), 4 (2007)
- [Ukk95] Ukkonen, E.: On-line construction of suffix trees. Algorithmica 14(3), 249–260 (1995)