

## 2. Preliminaries

**2.1. PROPERTIES OF SUFFIX TREES.** Let  $S \in \Sigma^n$  be a string formed by  $n$  characters drawn from  $\Sigma$ . The suffix tree  $T_S$  of  $S$  is the compacted trie<sup>2</sup> of all the suffixes of  $S\$, where  $\$ \notin \Sigma$ . Throughout the paper, we will assume that suffix trees are represented as follows. Leaf  $l_i$  represents suffix  $S[i, n]$ , and array entry  $l[i]$  points to  $l_i$ . Each internal node  $v$  has a length  $\sigma(v)$ , the sum of the edge lengths of its children. Then the length of  $v$ , denoted  $\sigma(v)$ , is  $\sigma(v) = \sum_{u \in v} \sigma(u) + 1$  where  $l_u$  is any leaf below  $v$ . The edges of nodes are stored in a list sorted by the first character on the edge from  $v$ . Neither these first characters, nor any other part of the string represented by any edge is stored, since it can be retrieved in linear time. See Figure 1 for an example. In this representation, the size of the suffix tree is  $O(n)$ .$

<sup>2</sup>A compacted trie differs from a trie in that maximal branch-free paths are replaced by edges labeled by the appropriate substring.

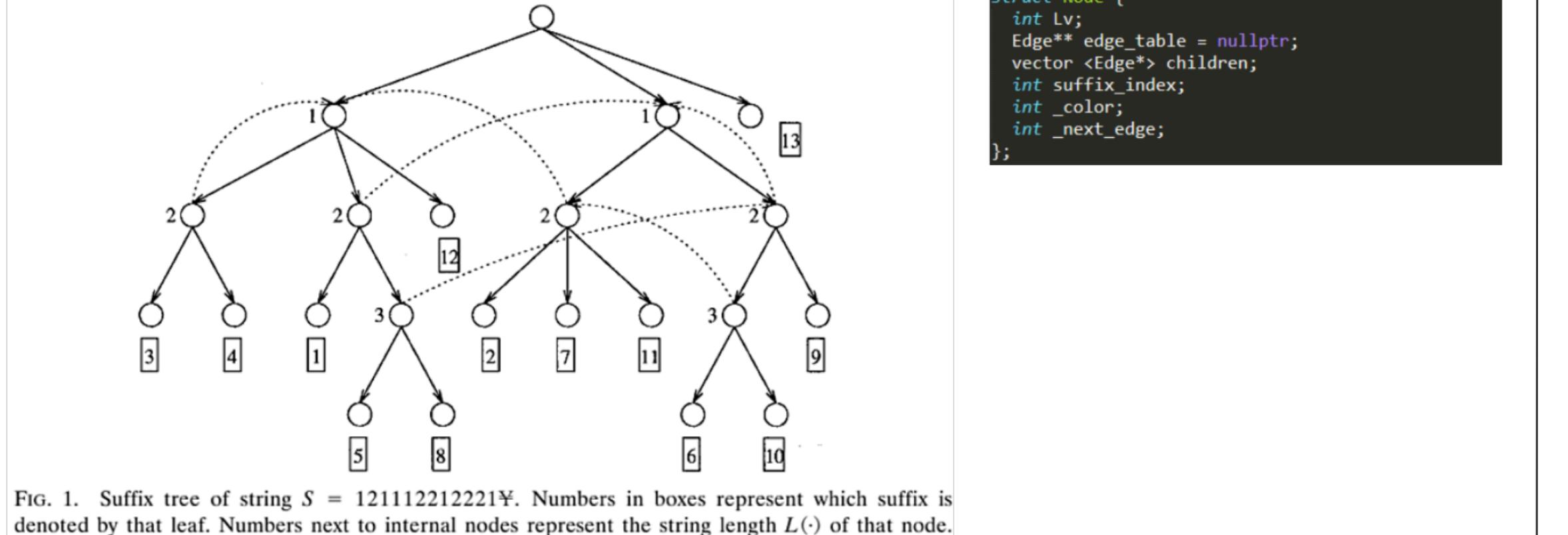


Fig. 1. Suffix tree of string  $S = 121112212221\$$ . Numbers in boxes represent which suffix is denoted by that leaf. Numbers next to internal nodes represent the string length  $L(v)$  of that node. Dotted lines are suffix links (see Lemma 2.1).

The following well-known lemma gives suffix trees a nice structure.

**LEMMA 2.1** [WEINER 1973]. Let  $a \in \Sigma$  and  $a \in \Sigma^*$ . If there is a node  $v$  in  $T_S$  such that  $l(v) = aa$ , then there is a node  $w$  in  $T_S$  such that  $l(w) = a$ .

Given this lemma we can define, for every node  $v$  in  $T_S$ , the suffix link  $l(v) = w$ , where  $v$  and  $w$  are defined as in Lemma 2.1. Notice that  $sl(v)$  links form a tree rooted at the root of  $T_S$ . The depth of any node  $v$  in this  $sl(v)$  tree is then just  $L(v)$ .

Given a string  $a \in \Sigma^n$ , let  $|a|$  denote the string length  $n$ . We use  $lcp(a, b)$  to indicate the length of the longest common prefix of any two strings  $a$  and  $b$ , and  $lca(v, w)$  to denote the least common ancestor of any two nodes  $v$  and  $w$  in a tree. The property of suffix trees most often exploited algorithmically is the following relationship between  $lcp$  in  $S$  and  $lca$  in  $T_S$ :

for all nodes  $v, w \in T_S$ ,  $lcp(l(v), l(w)) = |lca(v, w)|$ .

$lcp(l(v), l(w)) = \text{length of the longest common ancestor of } l(v) \text{ and } l(w)$

In all computational models, we will use the fact that it is easy to find the least common ancestors of two nodes [Harel and Tarjan 1984; Bender and Farach-Colton 2000; Chiang et al. 1995].

**2.2. SUFFIX ARRAYS AND EULER TOURS.** Let  $\Delta = \{S_i \mid S_i \in \Sigma^n\}$  be a set of strings of arbitrary lengths and assume, for the sake of exposition, that no string in  $\Delta$  is a prefix of another. We denote by  $lcp_\Delta$  the sorted compacted trie of the strings in  $\Delta$ . The *Euler Tour* of  $lcp_\Delta$  is a sequence of nodes in  $lcp_\Delta$  starting at the root and ending at the root. We denote by  $lcp_\Delta[i]$  the *i-th* node in  $lcp_\Delta$ . In particular, leaf  $i$  will now represent the string  $S_i$  and array entry  $l[i]$  will point to  $l_i$ . The size of  $T$  is  $O(|\Delta|)$  in addition to the overall length of the strings in  $\Delta$ . Since  $T$  is sorted, the children of any node  $v$  are stored in a list sorted by the first character on the edges from  $v$ . Observe that the in-order traversal of the leaves of  $T$  gives the lexicographically ordered sequence of the strings in  $\Delta$ . We denote by  $s_1, \dots, s_{|\Delta|}$  the permutation of  $1, \dots, |\Delta|$  such that  $S_{s_i}$  is lexicographically smaller than  $S_{s_j}$  if and only if  $s_i < s_j$ .

The suffix array of  $T$  consists of two arrays, the *sort array*  $\Lambda_T$  and the *longest common prefix array*  $LCP_T$ .  $\Lambda_T$  is the lexicographically ordered sequence of the strings in  $\Delta$ , that is  $\Lambda_T[i] = s_i$ . The array  $LCP_T$  stores the lengths of the common prefixes of adjacent strings in  $\Delta$ , that is  $LCP_T[i] = lcp(S_{s_i}, S_{s_{i+1}})$ .

Suppose we are given the sorted compacted trie  $T$  and we wish to compute  $\Lambda_T$  and  $LCP_T$ .  $\Lambda_T$  is simply the inorder listing of the leaves of  $T$ , whereas  $LCP_T$  can be obtained by means of Equality (1). Hence, both of them can be computed simultaneously during an in-order traversal of  $T$ . This is given by  $\Lambda_T$  and  $LCP_T$  if  $T$  is a straightforward search as recursive as the tree  $T$  in linear time in RAM [Bender and Muthukrishnan 1996]. We give other model-specific complexities below.

Given a rooted tree  $T' = (V, E)$ , the *Euler Tour*  $ET(T')$  of  $T'$  is a sequence of nodes,  $ET(T')$  has length  $2|E| + 1$  and is obtained by performing a depth-first search (DFS) on  $T'$  and outputting each node every time it is visited. In the RAM model,  $ET(T')$  is computed via an explicit DFS of the tree  $T'$ , whereas in the DAM model  $ET(T')$  is efficiently computed by simulating known PRAM-algorithms [Chiang et al. 1995].

## 3. RAM Algorithm

We start by developing a few tools that will simplify our presentation. Recall that  $lca(u, v)$  denotes the least common ancestor of a pair of nodes  $u, v$  in  $T$ . We will drop the  $T$ -subscript whenever there is no ambiguity.

**THEOREM 3.1** [HAREL AND TARJAN 1984; BENDER AND FARACH-COLTON 2000]. Given a sorted compacted trie  $T$  built on  $k$  strings,  $\Lambda_T$  and  $LCP_T$  can be preprocessed in  $O(nk)$  time so that, for any pair of nodes  $u, v$  in  $T$ , we have that  $lca(u, v) = lca(\Lambda_T[u], \Lambda_T[v])$ .

Given a sorted compacted trie  $T$  built on a set  $\Delta$  of  $k$  strings, we can compute the arrays  $\Lambda_T$  and  $LCP_T$  in  $O(k)$  time via a traversal of  $T$  (see Section 2.2). Conversely, given  $\Lambda_T$  and  $LCP_T$  we can reconstruct the sorted compacted trie  $T$  in linear time by exploiting Theorem 3.1 and Eq. (1). Notice that both transformations are independent of the length of the strings stored in the two data structures. Therefore, we can state:

**THEOREM 3.2.** Given a sorted compacted trie  $T$  built on  $k$  strings,  $\Lambda_T$  and  $LCP_T$  can be computed in  $O(k)$  time. The inverse transformation also takes  $O(k)$  time.

We have now all the ingredients to detail the implementation of the three steps of our divide-and-conquer approach for the case of the RAM model.

**3.1. BUILDING THE ODD TREE.** We show how to compute  $T_o$ , the compacted trie of all suffixes of  $S$  beginning in odd positions. Note that we can extend  $S$  with any character so that it has length a power of two. At the end of the computation, we simply prune the tree in linear time to remove the extra suffixes. The padding simplifies some boundary cases below, by making sure that the length of the string is always divisible by two. The padding at most doubles the length of  $S$  and so does not affect the complexity. While the string can be padded with any character, padding with  $\$$  makes it easier to strip away unwanted suffixes at the end of the computation.

*Step 1. Map pairs of characters into single characters as follows. For  $i = 1$  to  $n/2$ , pair  $(S[2i - 1], S[2i])$ . Lexicographically sort them by radix sort. May 26, 2022*

*Step 2. Map pairs of characters into single characters as follows. For  $i = 1$  to  $n/2$ , pair  $(S[2i - 1], S[2i])$ . Lexicographically sort them by radix sort. std::sort was used to perform a variation of bucket sort.*

*Step 3. Notice that any odd suffix  $S[2i - 1] \dots S[n/2]\$$  is equivalent to the suffix  $S'[i] \dots S[n/2]\$$ . Thus, by the lexicographic ordering of the characters of  $S$ , we have that  $\Lambda_{T_o}[i] = 2 \times \Lambda_T[i] - 1$ , and so this computation takes linear time. For our example,  $\Lambda_{T_o} = [3, 1, 5, 7, 11, 9, 13]$ .*

*Expectation is to reduce LCA to RMQ using a sparse table to answer LCP=LCA queries.*

*Step 3. Expand sorted suffix array.  $A_{T_o}[i] = 2 * A_{T_o}[i] - 1$*

*$A_{T_o} = [2, 1, 3, 4, 6, 5, 7]$*

```
/* expand lexigraphically ordered suffix array
void expand_sorted_suffix_array(vector<int>* A_ext, vector<int>* A) {
    for (int i = 0; i < A->size(); i++) {
        A_ext->push_back((A[i] + 1));
    }
}
```

*Step 4. Since each character in  $S'$  represents two characters in  $S$ , in order to compute  $LCP_{T_o}$  we need only to check whether the  $lcp$  between two adjacent suffixes in  $\Lambda_{T_o}$  must be extended by one unit. Thus,*

*$LCP_{T_o}[i] = 2 * LCP_T[i]$*

*+ 1 if  $S[\Lambda_{T_o}[i] + 2 * LCP_T[i]] = S[\Lambda_{T_o}[i + 1] + 2 * LCP_T[i]]$ ; 0 otherwise.*

*This computation takes linear time. For our example,  $LCP_{T_o}[i] = [1, 2, 0, 2, 1, 0]$ .*

*$LCP_{T_o}[0] = 2 * LCP_T[0] + 1$*

*$\{1, 0, 1, 0, 1, 0\}$*

*$LCP_{T_o}[1] = 2 * LCP_T[1] + 1$*

*$\{0, 1, 1, 0, 1, 0\}$*

*$LCP_{T_o}[2] = 2 * LCP_T[2] + 1$*

*$\{0, 0, 1, 1, 0, 0\}$*

*$LCP_{T_o}[3] = 2 * LCP_T[3] + 1$*

*$\{0, 0, 0, 1, 1, 0\}$*

*$LCP_{T_o}[4] = 2 * LCP_T[4] + 1$*

*$\{0, 0, 0, 0, 1, 1\}$*

*$LCP_{T_o}[5] = 2 * LCP_T[5] + 1$*

*$\{0, 0, 0, 0, 0, 1\}$*

*$LCP_{T_o}[6] = 2 * LCP_T[6] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[7] = 2 * LCP_T[7] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[8] = 2 * LCP_T[8] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[9] = 2 * LCP_T[9] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[10] = 2 * LCP_T[10] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[11] = 2 * LCP_T[11] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[12] = 2 * LCP_T[12] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[13] = 2 * LCP_T[13] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[14] = 2 * LCP_T[14] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[15] = 2 * LCP_T[15] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[16] = 2 * LCP_T[16] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[17] = 2 * LCP_T[17] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[18] = 2 * LCP_T[18] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[19] = 2 * LCP_T[19] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[20] = 2 * LCP_T[20] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[21] = 2 * LCP_T[21] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[22] = 2 * LCP_T[22] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[23] = 2 * LCP_T[23] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[24] = 2 * LCP_T[24] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[25] = 2 * LCP_T[25] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[26] = 2 * LCP_T[26] + 1$*

*$\{0, 0, 0, 0, 0, 0\}$*

*$LCP_{T_o}[27] = 2 * LCP_T[27] + 1$*