

Computing all repeats using suffix arrays[★]

František Franěk¹, W. F. Smyth^{1,2}, and Yudong Tang¹

¹ Algorithms Research Group, Department of Computing & Software
McMaster University, Hamilton, Ontario, Canada L8S 4K1
smyth@mcmaster.ca
www.cas.mcmaster.ca/cas/research/groups.shtml

² School of Computing, Curtin University, GPO Box U-1987
Perth WA 6845, Australia

June 28, 2002

Abstract. We describe an algorithm that identifies *all* the repeating substrings (tandem, overlapping, and split) in a given string $\mathbf{x} = \mathbf{x}[1..n]$. Given the suffix arrays of \mathbf{x} and of the reversed string $\hat{\mathbf{x}}$, the algorithm requires $\Theta(n)$ time for its execution and represents its output in $\Theta(n)$ space, either as a reduced suffix array (called an NE array) or as a reduced suffix tree (called an NE tree). The output substrings \mathbf{u} are *nonextendible* (NE); that is, extension of some occurrence of \mathbf{u} in \mathbf{x} , either to the left or to the right, yields a string ($\lambda\mathbf{u}$ or $\mathbf{u}\lambda$) that is unequal to the same extension of some other occurrence of \mathbf{u} . Thus the number of substrings output is the minimum required to identify all the repeating substrings in \mathbf{x} . The output can be used in a straightforward way to identify only repeating substrings that satisfy some proximity or minimum length condition.

1 Introduction

The computation of all the repeating substrings in a given string $\mathbf{x} = \mathbf{x}[1..n]$ is a problem with various application areas, most notably data compression, cryptography, and computational biology. For repeating substrings that are *tandem* (that is, *repetitions*), several $O(n \log n)$ algorithms [1, 5, 10] were discovered about 20 years ago; more recently, a repetitions algorithm [9] was published that, at least theoretically, executes in $\Theta(n)$ time in the common case that

[★] Supported in part by grants from the Natural Sciences & Engineering Research Council of Canada.

the alphabet is *indexed* — that is, treatable as a range of integers $1..n$.

These successes with repetitions have encouraged researchers to seek algorithms that efficiently compute all repeating substrings \mathbf{u} , including in addition to tandem occurrences those that are *split* (of the form \mathbf{uvu} for some nonempty \mathbf{v}) and *overlapping* (such as $\mathbf{u} = \mathbf{abaab}$ in $\mathbf{x} = \mathbf{abaabaab}$). The following definitions permit this problem to be stated more precisely:

Definition 1. A *repeat* in a string \mathbf{x} is a tuple

$$M_{\mathbf{x},\mathbf{u}} = (p; i_1, i_2, \dots, i_r), \quad r \geq 2,$$

where

$$\mathbf{u} = \mathbf{x}[i_1..i_1+p-1] = \mathbf{x}[i_2..i_2+p-1] = \dots = \mathbf{x}[i_r..i_r+p-1].$$

\mathbf{u} is said to be a *repeating substring* of \mathbf{x} and the *generator* of $M_{\mathbf{x},\mathbf{u}}$. As for repetitions, we call $p = |\mathbf{u}|$ the *period* of the repeat and r its *exponent*. If the tuple includes all the occurrences of \mathbf{u} in \mathbf{x} , then $M_{\mathbf{x},\mathbf{u}}$ is said to be *complete* and is written $M_{\mathbf{x},\mathbf{u}}^*$.

In its most general form, then, our task may be defined to be the computation of $M_{\mathbf{x},\mathbf{u}}^*$ for every repeating \mathbf{u} . This task is simplified by introducing the idea of “nonextendibility”:

Definition 2. A repeat $M_{\mathbf{x},\mathbf{u}} = (p; i_1, i_2, \dots, i_r)$ is said to be *left-extendible* (respectively, *right-extendible*) if and only if

$$(p+1; i_1-1, i_2-1, \dots, i_r-1) \quad (\text{respectively, } (p+1; i_1, i_2, \dots, i_r))$$

is a repeat. If $M_{\mathbf{x},\mathbf{u}}$ is neither left-extendible nor right-extendible, it is said to be *nonextendible*. We abbreviate these terms as **LE**, **RE** and **NE**, respectively.

It is clear that it suffices to compute complete NE repeats: those that are extendible, either to left or right, do not need to be reported because any extendible repeating substring will be reported implicitly as a substring of the strings specified by an NE repeat.

In the last few years, two algorithms [8, 3] have been published that employ suffix trees to compute all the NE repeats in a given

string \mathbf{x} defined on an ordered alphabet. There are several standard $O(n \log n)$ -time algorithms [12–14] that compute suffix trees: indeed, if alphabet size is assumed to be constant, then the $\log n$ term can be defined away and the claim made that these algorithms are “linear” in string length. It is not usually recognized that Crochemore’s repetitions algorithm [5] also belongs in this collection of $O(n \log n)$ -time algorithms: by continuing the refinement process until all of the equivalence classes become singletons, all the information required for suffix tree construction becomes available. For strings defined on an indexed alphabet, there does exist a genuine $\Theta(n)$ -time suffix tree construction algorithm [7], but it is very complicated and the space requirements are consequently so great that its theoretical asymptotic time advantage is eaten away by additional memory swaps to and from disk. In general, the main difficulty with suffix tree use is the space requirement, many times that of the string itself, so that for strings of several tens of millions of characters, their use becomes impractical. As a result, related data structures such as directed acyclic word graphs (DAWGs) [4, 6] and suffix arrays [11] have been devised, that provide much of the functionality of suffix trees but use less space.

In its simplest form, a suffix array $\sigma_{\mathbf{x}} = \sigma[1..n]$ of a given string $\mathbf{x}[1..n]$ is an array of integers such that for every $i \in 1..n$, $\sigma[i] = j$ if and only if $\mathbf{x}[j..n]$ is the i^{th} smallest suffix of $\mathbf{x}\$$, where the sentinel letter $\$$ is by convention larger than any other letter in the alphabet. Thus $\sigma_{\mathbf{x}}$ specifies the starting positions of the suffixes of \mathbf{x} in ascending lexicographical order; for example, if

$$\begin{array}{cccccccccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 \\ \mathbf{x} & = & a & b & a & a & b & a & b & a & a & b & a & a & b, \end{array} \quad (1)$$

then

$$\sigma_{\mathbf{x}} = 8\ 3\ 11\ 6\ 1\ 9\ 4\ 12\ 7\ 2\ 10\ 5\ 13. \quad (2)$$

For effective use of suffix arrays, a second array $\lambda_{\mathbf{x}}$ is normally also computed that gives the length (lcp) of the longest common prefix (LCP) of adjacent entries in $\sigma_{\mathbf{x}}$; in our example,

$$\lambda_{\mathbf{x}} = 0431653205421, \quad (3)$$

reflecting the fact that $\text{lcp}(8, 3) = 4$, $\text{lcp}(3, 11) = 3$, and so on. In general, we set $\lambda[1] = 0$, while for every $i \in 2..n$,

$$\lambda[i] = \text{lcp}(\sigma[i-1], \sigma[i]). \quad (4)$$

Suffix arrays were introduced in [11], where an $O(n \log n)$ -time and space-efficient algorithm is described for their calculation. For many problems — in particular, as we shall see, for the problem discussed in this paper — suffix arrays provide the same functionality that suffix trees make available, and with a very much reduced use of space. The drawback is that direct suffix array construction requires, according to [11], three to ten times the time needed for construction of the corresponding suffix tree. Of course suffix arrays can be computed from suffix trees, but then the space advantage is lost.

In this paper we describe an alternate approach to the calculation of all the NE repeats in a string \mathbf{x} that can be used either with suffix trees or suffix arrays, depending on circumstances. In fact, a weaker form of suffix array/tree (which we call a *quasi suffix array/tree*) suffices for this algorithm: there is no requirement for the suffixes of \mathbf{x} to be sorted in lexicographic order, and so a simpler and faster suffix array construction algorithm, such as that proposed in [2], can be used. Referring again to the example (1), this means that we might use

$$\begin{aligned} \sigma_{\mathbf{x}} &= 1\ 6\ 9\ 4\ 12\ 3\ 8\ 11\ 2\ 7\ 10\ 5\ 13 \\ \lambda_{\mathbf{x}} &= 0\ 6\ 5\ 3\ 2\ 1\ 4\ 3\ 0\ 5\ 4\ 2\ 1 \end{aligned} \quad (5)$$

just as well as the arrays (2) and (3).

The new algorithm is based on the following easily-proved lemma:

Lemma 1. *Let $\hat{\mathbf{x}}$ denote the reverse string $\mathbf{x}[n]\mathbf{x}[n-1] \cdots \mathbf{x}[1]$ of a given string \mathbf{x} . Then a repeat $M_{\mathbf{x}, \mathbf{u}}$ is LE if and only if $M_{\hat{\mathbf{x}}, \hat{\mathbf{u}}}$ is RE. \square*

This result suggests a straightforward approach to computing all the NE repeats in \mathbf{x} :

- compute all the NRE (non-RE) repeats of \mathbf{x} and all the NRE (non-RE) repeats of $\hat{\mathbf{x}}$;

- compare the NRE repeats of \mathbf{x} with the NRE repeats of $\widehat{\mathbf{x}}$ to identify those that are in both lists (the NE repeats).

It turns out that the NRE repeats of a string are identified by its suffix tree (as well as by its suffix array), a fact illustrated by reference to our example (1) in Figure 1. This figure represents a kind of stripped-down suffix tree in which we represent only the distinct prefixes of each suffix in \mathbf{x} . The internal (circular) nodes are recognizable as the lcp values contained in the array $\lambda_{\mathbf{x}}$, while the leaf (square) nodes represent the starting positions i of the suffixes $\mathbf{x}[i..n]$, $i = 1, 2, \dots, n$, as specified in the array $\sigma_{\mathbf{x}}$. As in every suffix tree, the lcp value at the root of each subtree applies to all the leaf nodes of that subtree; thus the lcp values also specify the complete repeats in \mathbf{x} that are according to Definition 2 NRE. For example, the substrings

$$\mathbf{u} = \mathbf{x}[9..13] = \mathbf{x}[6..10] = \mathbf{x}[1..5]$$

defined by the subtree rooted at the internal node 5 constitute a complete NRE repeat

$$M_{\mathbf{x},\mathbf{u}}^* = (5; 9, 6, 1).$$

In Section 2 we describe an algorithm that computes the NE tree of \mathbf{x} , hence all the NE repeats in \mathbf{x} , in linear time from the suffix (NRE) trees of \mathbf{x} and $\widehat{\mathbf{x}}$. Then Section 3 explains how this algorithm can be implemented to compute the NE array of \mathbf{x} , also in linear time, based on the suffix (NRE) arrays of \mathbf{x} and $\widehat{\mathbf{x}}$. Section 4 briefly discusses extensions of these algorithms and open problems.

2 Computing the NE Tree

We suppose that the suffix trees, say $T_{\mathbf{x}}$ and $T_{\widehat{\mathbf{x}}}$ of \mathbf{x} and $\widehat{\mathbf{x}}$ respectively, have already been computed, thereby identifying complete NRE repeats in \mathbf{x} and complete NRE repeats in $\widehat{\mathbf{x}}$. From Lemma 1 we know that the NRE repeats in $\widehat{\mathbf{x}}$ identify NLE repeats in \mathbf{x} . For brevity we refer to the generator \mathbf{u} of an NRE (respectively, NLE, NE) repeat $M_{\mathbf{x},\mathbf{u}}$ as an NRE (respectively, NLE, NE) repeating substring, while asking the reader to bear in mind that the terms NRE,

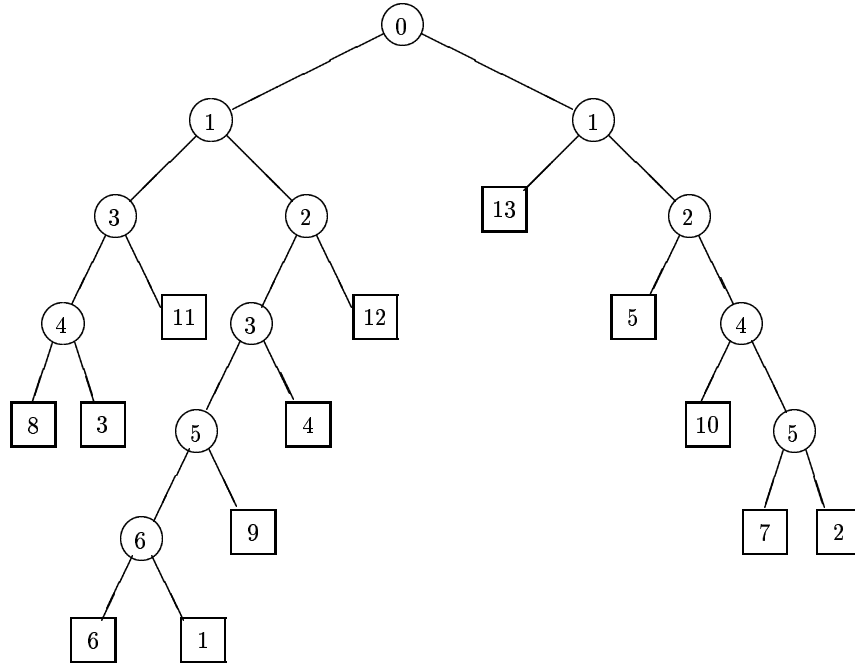


Fig. 1. Suffix (NRE) tree $T_{\mathbf{x}}$ for $\mathbf{x} = abaababaabaab$

NLE, NE really relate to the repeat that is a collection of repeating substrings.

Suppose that some lcp node \hat{p} in $T_{\hat{\mathbf{x}}}$ has as child a position node \hat{i} . Then in $\hat{\mathbf{x}}$ there exists an NRE repeating substring

$$\hat{\mathbf{u}} = \hat{\mathbf{x}}[\hat{i}..\hat{i}+\hat{p}-1]$$

that is the reverse of an NLE repeating substring

$$\mathbf{u} = \mathbf{x}[n-(\hat{i}+\hat{p}-1)+1..n-\hat{i}+1]$$

in \mathbf{x} . Thus the assignment

$$i \leftarrow n - (\hat{i} + \hat{p} - 2) \tag{6}$$

identifies one start position i of an NLE repeating substring \mathbf{u} in \mathbf{x} of period \hat{p} . Observe that in fact i is also a start position of one of a set of NLE repeating substrings of every period $j \in 1..\hat{p}$; thus if in $T_{\mathbf{x}}$ we find for some i the largest value of \hat{p} such that i is one of a collection

of substrings of length \widehat{p} that are both NRE and NLE, then every parent of \widehat{p} in $T_{\mathbf{x}}$ will also identify collections of substrings that are both NRE and NLE.

Algorithm 1 (Compute the NE Tree)

— Given the suffix trees $T_{\mathbf{x}}$ & $T_{\widehat{\mathbf{x}}}$, compute the NE tree of $\mathbf{x}[1..n]$

- (1) traverse $T_{\mathbf{x}}$ to create a table `POINTER`[i] that for each position i in \mathbf{x} points to the corresponding node in $T_{\mathbf{x}}$
- (2) **for** every lcp node p in $T_{\mathbf{x}}$ **do**
`NE`[p] \leftarrow `FALSE`
- (3) **for** every parent-child pair $(\widehat{p}, \widehat{i})$ in $T_{\widehat{\mathbf{x}}}$ **do**
 — Here use `POINTER`[i]:
if $i = n - (\widehat{i} + \widehat{p} - 2)$ is a child of lcp node \widehat{p} in $T_{\mathbf{x}}$ **then**
 while not `NE`[\widehat{p}] **do**
 `NE`[\widehat{p}] \leftarrow `TRUE`
 if $\widehat{p} \neq 0$ **then**
 $\widehat{p} \leftarrow$ parent of \widehat{p} in $T_{\mathbf{x}}$
- (4) traverse $T_{\mathbf{x}}$ deleting every subtree rooted at an lcp node p for which `NE`[p] = `FALSE`

Let us say that a substring \mathbf{u} of \mathbf{x} is a *maximal NE repeating substring* of \mathbf{x} if and only if

- \mathbf{u} occurs at least twice in \mathbf{x} ;
- \mathbf{u} is not a proper substring of any repeating substring of \mathbf{x} .

The following result, again easily proved, then tells us that a maximal NE repeating substring must be identifiable in both $T_{\mathbf{x}}$ and $T_{\widehat{\mathbf{x}}}$:

Lemma 2. *If $\mathbf{u} = \mathbf{x}[i..i+p-1]$ is a maximal NE repeating substring of \mathbf{x} , then the position node i occurs as a child of the lcp node p in $T_{\mathbf{x}}$, and $n - (i - p - 2)$ occurs as a child of p in $T_{\widehat{\mathbf{x}}}$. \square*

Recall the observation made above that every internal (lcp) node in $T_{\mathbf{x}}$ that is an ancestor of an NE repeating substring must itself be the root of a subtree of $T_{\mathbf{x}}$ whose leaf (position) nodes are the repeating substrings of a complete NE repeat. This lemma therefore provides us with a simple strategy for identifying all NE repeating substrings: just find the maximal ones (largest value of p), then locate all their

ancestors in $T_{\mathbf{x}}$. Algorithm 1 provides an outline of how this can be accomplished.

The algorithm is expressed in terms of four steps. Step (1) is a traversal of $T_{\mathbf{x}}$ that sets up a table enabling each position node i in $T_{\mathbf{x}}$ to be accessed later in constant time. In Step (2) a Boolean variable corresponding to each mark node p in $T_{\mathbf{x}}$ is initialized to **FALSE**, indicating that no terminal nodes in the subtree rooted at p have currently been identified as NE. Step (3) processes every parent-child pair (\hat{p}, \hat{i}) in $T_{\hat{\mathbf{x}}}$, testing to determine whether or not the equivalent parent-child pair $(\hat{p}, n - (\hat{i} + \hat{p} - 2))$ exists in $T_{\mathbf{x}}$; if so, then the mark node \hat{p} and all its ancestors in $T_{\mathbf{x}}$ must be NE — accordingly, until an ancestor is found that is already NE, \hat{p} and its ancestors in $T_{\mathbf{x}}$ are identified as NE. A final step traverses $T_{\mathbf{x}}$ to eliminate all subtrees rooted at any node p for which $\text{NE}[p] = \text{FALSE}$; the remaining tree, $T_{\mathbf{x}}^{\text{NE}}$, is the NE tree of \mathbf{x} .

Each of the steps (1), (2), (4) of Algorithm 1 is a traversal of an NRE tree that requires $O(1)$ time at each node, hence $O(n)$ time in total. Step (3) is another tree traversal, but with slightly more complex processing: since the **if** statement makes use of the **POINTER** array, it too requires only $O(1)$ time at each execution, hence also $O(n)$ time overall. The **while** loop in Step (3) causes the NE variable for at most n mark nodes to be set **TRUE**, and tests the current setting of at most n mark nodes for which NE is already **TRUE** — thus the total time consumed in the while loop is $O(n)$. We have established:

Theorem 1. *Given the suffix trees $T_{\mathbf{x}}$ and $T_{\hat{\mathbf{x}}}$ for $\mathbf{x}[1..n]$, Algorithm 1 correctly computes the NE tree $T_{\mathbf{x}}^{\text{NE}}$ using $O(n)$ time and $\Theta(n)$ additional space. \square*

Figure 2 displays the suffix tree $T_{\hat{\mathbf{x}}}$ for our example string (1), modified so that position nodes \hat{i} are replaced by $i = n - (\hat{i} + \hat{p} - 2)$. Note that only positions 1, 6, 9 give rise to NLE repeats; thus the NE tree reduces to the subtree of $T_{\mathbf{x}}$ that corresponds to these positions — that is, the path from the root of $T_{\mathbf{x}}$ to lcp node 6 that corresponds to the substring *abaaba*. The complete NE repeats of (1) are thus seen from Figure 1 to be all the occurrences of *a*, *ab*, *aba*, *abaab* and *abaaba*. In general, the occurrences of NE repeats, whether all of them or selected according to various criteria, can be located and

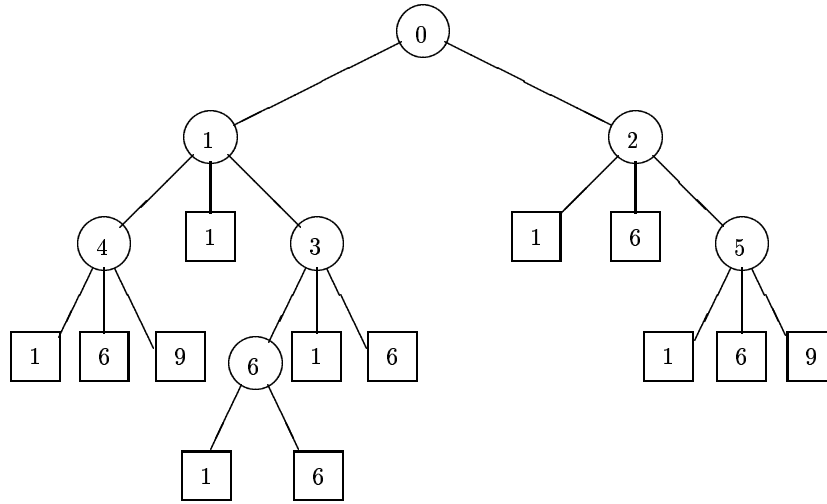


Fig. 2. $T_{\hat{x}}$ for $\hat{x} = baabaababaaba$ with $i = n - (\hat{i} + \hat{p} - 2)$

output in time proportional to their number by traversal of the NE tree. Such problems are also discussed in [3]. Another quite different approach to the computation of nonextendible repeats (though still making use of suffix trees) is described in [8].

3 Computing the NE Array

We provide here an outline of a version of the algorithm that avoids the use of suffix trees completely: the entire calculation is carried out based on suffix arrays (both σ_x and λ_x), and its end-product is not an NE tree, but rather an NE array. Since the suffix array of $x[1..n]$ requires only $2n$ computer words for its storage, this revision yields significant saving in space.

For every $j \in 1..n$, let $p_j = \lambda[j]$. It then follows from (4) that for every $j \in 2..n$,

```

if  $p_j > p_{j-1}$  then
     $p_j$  is a descendant of  $p_{j-1}$  in the suffix tree
elseif  $p_j < p_{j-1}$  then
     $p_j$  is an ancestor of  $p_{j-1}$  in the suffix tree
else
     $p_j$  and  $p_{j-1}$  identify the same node in the suffix tree

```

Since the suffix array contains the same lcp information that is provided by the suffix tree, these relationships imply that the suffix array must be decomposable into subarrays, each of which specifies nodes that lie on a single path from the root to a terminal node in the suffix tree. For example, the lcp values in the array (3) can be separated into three subarrays

$$0, 4, 3, \underline{1} / \underline{1}, 6, 5, 3, 2, \underline{0} / \underline{0}, 5, 4, 2, 1$$

corresponding to the three paths in the suffix tree of Figure 1. Here the repeated lcp values $\underline{0}$ and $\underline{1}$ identify the root of a subtree at which a new path begins. These nodes are called **branch nodes**, marking the point at which a new path diverges from the preceding one. A similar breakdown is provided by the lcp array (5):

$$0, 6, 5, 3, 2, \underline{1} / \underline{1}, 4, 3, \underline{0} / \underline{0}, 5, 4, 2, 1.$$

These ideas can be made more precise as follows:

Definition 3. An **I-run** $I_{j,h}$ in a suffix array is a sequence of $h \geq 2$ lcp values $p_j, p_{j+1}, \dots, p_{j+h-1}$, not all equal, such that

- (a) either $j = 1$ or else $p_{j-1} > p_j$;
- (b) $p_j \leq p_{j+1} \leq \dots \leq p_{j+h-1}$;
- (c) either $j+h-1 = n$ or else $p_{j+h-1} > p_{j+h}$.

Thus an I-run is a maximum-length sequence of nondecreasing lcp values in the suffix array; similarly, a **D-run** $D_{j,h}$ is a maximum-length sequence of nonincreasing lcp values. Using these definitions, we can now characterize paths in the suffix tree in terms of runs in the suffix array:

Lemma 3. Every sequence $I_{j,h}D_{j+h-1,h'}$ or $I_{j,h}$, where $j+h-1 = n$, in the suffix array identifies nodes that lie on a single path in the suffix tree.

Proof. Since $p_1 = 0$, it follows from Definition 3(a)-(b) that either there are no runs in the suffix array ($p_j = 0$ for every $j \in 1..n$) or else the sequence of runs begins with an I-run $I_{1,h}$. Observe that Definition 3(c) ensures that every I-run is followed by a D-run or else by the end of the array; similarly, every D-run must be followed by an I-run or else by the end of the array.

Now consider an I-run $I_{j,h}$. Since $p_{j+(t-1)} \leq p_{j+t}$ for every $t \in 1..h-1$, it follows that $I_{j,h}$ determines a sequence of descendants of p_j that lie on a single path in the suffix tree. If in fact $i+h-1 = n$, then $I_{j,h}$ is in itself the final path identified in the suffix array. If not, however, then $I_{j,h}$ is followed by a D-run $D_{j+h-1,h'}$, where $p_{(j+h-1)+(t-1)} \leq p_{(j+h-1)+t}$ for every $t \in 1..h'-1$. These lcp values determine a sequence of ancestors of p_{j+h-1} , all on the same path from the root already identified by $I_{j,h}$. The subsequent I-run $I_{j+h+h'-2,h''}$, if it exists, of course specifies descendants of the branch node $p_{j+h+h'-2}$ that lie on a path distinct from the path specified by $I_{j,h}$. \square

We see then that the ID-runs in the suffix array essentially decompose the internal (lcp) nodes of the suffix array into paths. Adjacent pairs of runs in the array identify two paths in the tree that have exactly one node (a branch node) in common. Observe also that for each ID-run, the minimum node in the corresponding path must be either the first node of the I-run or the last node of the D-run.

Since a branch node is *both* the last node of a D-run *and* the first node of an I-run, the values of pairs of adjacent branch nodes in the suffix array therefore determine the way in which two paths in the suffix tree are related. If we include the first lcp value ($p_1 = 0$) in the suffix array as a branch node, and denote consecutive branch nodes by b_1 and b_2 , respectively, we see that there are only three possible relationships between the corresponding paths, as shown in Figure 3. Notice in particular that in each of these three cases, future paths can be added only as descendants of b_2 on the current path or as descendants of ancestors of b_2 in the tree. No further change can be made to the path rooted at b_2 that was specified by the previous path containing b_1 .

This discussion puts us in a position to describe an approach to the calculation of NE repeats that is based on suffix arrays only. As before we denote the position and lcp arrays for \mathbf{x} by $\sigma_{\mathbf{x}}$ and $\lambda_{\mathbf{x}}$, respectively; for $\widehat{\mathbf{x}}$ the corresponding arrays are $\widehat{\sigma}_{\mathbf{x}}$ and $\widehat{\lambda}_{\mathbf{x}}$. We also make use of a LOC array that gives for each position node \widehat{i} in $T_{\widehat{\mathbf{x}}}$ its location in $\widehat{\sigma}_{\mathbf{x}}$; this enables us to find any position in $\widehat{\sigma}_{\mathbf{x}}$ in constant time. A bit vector $\text{NE}[1..n]$ specifies for each $j \in 1..n$ whether (TRUE) or not (FALSE) the lcp node $\lambda[j]$ is also NLE. Finally, in order to

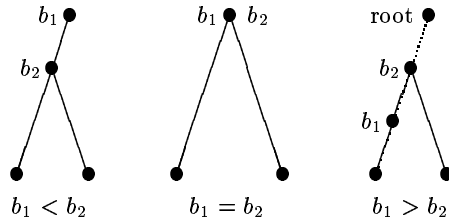


Fig. 3. Branch nodes b_1 and b_2 of adjacent paths (ID-runs)

identify the beginning and end of ID paths in $\lambda_{\mathbf{x}}$, we use the ID patterns (Lemma 3) to compute a BRANCH array that includes, in addition to the branch points defined above, also the first and last positions in $\lambda_{\mathbf{x}}$. An outline of the processing, more complex than the tree-based calculation but not necessarily more time-consuming, is displayed as Algorithm 2.

Step (1) of the algorithm is a straightforward **for** loop that initializes the LOC array and the bit vector NE. Step (2) computes the BRANCH array based on Lemma 3 and a linear scan of $\lambda[1..n]$. Step (3) is another **for** loop that first identifies all the matches of (p, i) pairs in the NRE arrays for \mathbf{x} and $\hat{\mathbf{x}}$, then for each path ensures that every ancestor of a matched (p, i) pair within the path is NE. Each of the two sections of Step (3) requires $\Theta(n)$ time. In Step (4) the processing of the second section of Step (3) is repeated, but in reverse order to ensure that all ancestors of NE positions are properly set to be themselves NE. Finally the NE array is determined by all those positions j in $\lambda_{\mathbf{x}}$ for which $\text{NE}[j] = \text{TRUE}$. Hence

Theorem 2. *Given the suffix (NRE) arrays corresponding to $\mathbf{x}[1..n]$ and $\hat{\mathbf{x}}$, Algorithm 2 correctly computes the NE array using $\Theta(n)$ time and $\Theta(n)$ additional space. \square*

We remark that the NE array can if desired be used to compute the corresponding NE tree in linear time.

4 Extensions & Open Problems

In this paper we have described an algorithm — implemented in two quite different ways — that, given suffix trees or suffix arrays of the

Algorithm 2 (Compute the NE Array)

— Given the suffix (NRE) arrays of $\mathbf{x} = \mathbf{x}[1..n]$ and $\hat{\mathbf{x}}$, compute the NE array of \mathbf{x}

- (1) **for** $j \leftarrow 1$ **to** n **do**
 $\text{LOC}[\hat{\sigma}[j]] \leftarrow j$; $\text{NE}[j] \leftarrow \text{FALSE}$
- (2) compute $\text{BRANCH}[1..b^*]$, where for $b \in 1..b^*$, $\text{BRANCH}[b]$ is the position of the b^{th} branch point in $\lambda_{\mathbf{x}}$ (include positions $j = 1$ and $j = n$ as branch points)
- (3) $b \leftarrow 1$
 for $j \leftarrow 1$ **to** n **do**
 determine j' such that $(p, i) = (\lambda[j'], \sigma[j])$ is a parent-child pair in the suffix tree of \mathbf{x}
 $\hat{i} \leftarrow n - (i + p - 2)$; $\hat{j} \leftarrow \text{LOC}[\hat{i}]$
 determine \hat{j}' such that $(\hat{p}, \hat{i}) = (\hat{\lambda}[\hat{j}'], \hat{\sigma}[\hat{j}])$ is a parent-child pair in the suffix tree of $\hat{\mathbf{x}}$
 if $p = \hat{p}$ **then**
 $\text{NE}[j'] \leftarrow \text{TRUE}$; $\text{NE}[1] \leftarrow \text{TRUE}$
 if $j = \text{BRANCH}[b+1]$ **then**
 determine $j^* \in \text{BRANCH}[b].. \text{BRANCH}[b+1]$ such that
 $p^* = \lambda[j^*]$ is the maximum lcp
 for which $\text{NE}[j^*] = \text{TRUE}$
 for every $j' \in \text{BRANCH}[b].. \text{BRANCH}[b+1]$ **do**
 if $\lambda[j'] \leq p^*$ **then**
 $\text{NE}[j'] \leftarrow \text{TRUE}$
 $b \leftarrow b+1$
- (4) process the paths identified by BRANCH in reverse order, setting NE true for ancestors of p^* as in Step (3)

strings \mathbf{x} and $\widehat{\mathbf{x}}$, computes all the complete NE repeats in linear time while representing them in linear space. To restrict the output to repeats that lie within a specified substring of \mathbf{x} , it is necessary only to ignore values of $\sigma[i]$ (position values) that lie outside the confines of that substring; similarly, in order to consider only repeats above a certain length ℓ , it suffices to consider only subtrees whose roots contain lcp values $\lambda[i]$ greater than ℓ .

Although we have described the algorithms in terms of suffix arrays and suffix trees, we remark again that quasi suffix arrays (and of course corresponding quasi suffix trees) such as the one illustrated in (5) can be used instead to implement Algorithms 1 and 2. This means that more efficient algorithms such as the one described in [2] can be used.

We remark also that since the order of output is determined by the structure of the suffix tree/array, it is therefore not straightforward to output repeating substrings in the natural order of their left-to-right occurrence in \mathbf{x} . We leave this as an open problem.

References

- [1] Alberto Apostolico & Franco P. Preparata, **Optimal off-line detection of repetitions in a string**, *TCS 22* (1983) 297-315.
- [2] Leila Baghdadi, František Franěk, W. F. Smyth & Xiangdong Xiao, **Computing quasi suffix arrays**, preprint (2002).
- [3] Gerth S. Brodal, Rune B. Lyngso, Christian N. S. Pedersen & Jens Stoye, **Finding maximal pairs with bounded gap**, *J. Discrete Algs. 1* (2000) 77-103.
- [4] M. T. Chen & Joel Seiferas, **Efficient and elegant subword-tree construction**, *Combinatorial Algorithms on Words (NATO ASI Series F12)*, Alberto Apostolico & Zvi Galil (eds.), Springer-Verlag (1985) 97-107.
- [5] Maxime Crochemore, **An optimal algorithm for computing the repetitions in a word**, *IPL 12-5* (1981) 244-250.
- [6] Maxime Crochemore & R. Verin, **Direct construction of compact directed acyclic word graphs**, *Proc. Eighth Annual Symp. Combinatorial Pattern Matching*, Lecture Notes in Computer Science 1264, Springer-Verlag (1997) 116-129.
- [7] Martin Farach, **Optimal suffix tree construction with large alphabets**, *Proc. 38th Annual IEEE Symp. Foundations of Computer Science* (1997) 137-143.
- [8] Dan Gusfield, *Algorithms on Strings, Trees, & Sequences*, Cambridge University Press (1997) 534 pp.
- [9] Roman Kolpakov & Gregory Kucherov, **On maximal repetitions in words**, *J. Discrete Algs. 1* (2000) 159-186.
- [10] Michael G. Main & Richard J. Lorentz, **An $O(n \log n)$ algorithm for finding all repetitions in a string**, *J. Algs. 5* (1984) 422-432.
- [11] Udi Manber & Gene W. Myers, **Suffix arrays: a new method for on-line string searches**, *SIAM J. Comput. 22-5* (1993) 935-948.

- [12] Edward M. McCreight, **A space-economical suffix tree construction algorithm**, *JACM* 32-2 (1976) 262-272.
- [13] Esko Ukkonen, **Constructing suffix trees on-line in linear time**, *Proc. IFIP 92*, vol. I (1992) 484-492.
- [14] P. Weiner, **Linear pattern matching algorithms**, *Proc. 14th Annual IEEE Symp. Switching & Automata Theory* (1973) 1-11.