

International Journal of Foundations of Computer Science
© World Scientific Publishing Company

Reachability on Suffix Tree Graphs

Yasuto Higa

*Department of Informatics, Kyushu University, 744 Motoooka
Nishi-ku Fukuoka 819-0395 JAPAN
y-higa@i.kyushu-u.ac.jp
<http://http://www.i.kyushu-u.ac.jp/~ie204013/>*

Hideo Bannai*

*Department of Informatics, Kyushu University, 744 Motoooka
Nishi-ku Fukuoka 819-0395 JAPAN
bannai@i.kyushu-u.ac.jp*

Shunsuke Inenaga†

*Graduate School of Information Science and Electrical Engineering, Kyushu University,
Nishi-ku Fukuoka 819-0395 JAPAN
shunsuke.inenaga@i.kyushu-u.ac.jp*

Masayuki Takeda‡

*Department of Informatics, Kyushu University, 744 Motoooka
Nishi-ku Fukuoka 819-0395 JAPAN
takeda@i.kyushu-u.ac.jp*

Received (Day Month Year)
Accepted (Day Month Year)
Communicated by (xxxxxxxxxx)

We analyze the complexity of graph reachability queries on *ST-graphs*, defined as directed acyclic graphs (DAGs) obtained by merging the suffix tree of a given string and its suffix links. Using a simplified reachability labeling algorithm presented by Agrawal *et al.* (1989), we show that for a random string of length n , its *ST-graph* can be preprocessed in $O(n \log n)$ expected time and space to answer reachability queries in $O(\log n)$ time. Furthermore, we present a series of strings that require $\Theta(n\sqrt{n})$ time and space to answer reachability queries in $O(\log n)$ time for the same algorithm. Exhaustive computational calculations for strings of length $n \leq 33$ have revealed that the same strings are also the worst case instances of the algorithm. We therefore conjecture that reachability queries can be answered in $O(\log n)$ time with a worst case time and space preprocessing

*Partly supported by The Ministry of Education, Culture, Sports, Science and Technology, Grant-in-Aid for Young Scientists (B)

†Partly supported by JSPS Research Fellowships for Young Scientists.

‡SORST, Japan Science and Technology Agency (JST)

2 Yasuto Higa et al.

complexity of $\Theta(n\sqrt{n})$.

Keywords: algorithms and data structures, suffix trees, graph reachability.

1. Introduction

The reachability query for two nodes u, v of a given directed graph is to answer whether or not there exists a path in the graph that starts from u and ends at v . For any given graph, the query can be answered in $O(n + e)$ time by conducting a simple depth-first traversal on the graph, where n is the number of nodes and e is the number of edges in the graph.

There have been several studies on preprocessing a graph in order to answer reachability queries more efficiently [5, 1, 7, 3, 9]. A simple approach would be to construct the transitive closure of the graph, achieving $O(1)$ time query at the cost of $O(n^2)$ time and space for the preprocessing. For graphs with specific properties, there exists methods with smaller complexity bounds. Graph reachability for planar graphs with a single source node and sink node was considered in [5]. Reachability queries for such graphs can be answered in $O(1)$ time given $O(n + e)$ time and $O(n)$ space preprocessing. For partial lattices, a method which achieves $O(1)$ time query with $O(n^2)$ time and $O(n\sqrt{n})$ space preprocessing was shown in [7], where n is the size of the ground set. When considering arbitrary graphs with n vertices and e edges, it has been shown in [3] that for any labeling scheme, there exists a graph such that the reachability labeling has total size of $\Omega(n\sqrt{e})$.

In this paper, we consider the graph reachability query problem on *ST-graphs*, which are DAGs derived from suffix trees and suffix links. *ST-graphs* are neither planar nor partial lattices. A suffix tree of a given string is a data structure that captures important information concerning the substrings of the string [10]. We present and analyze a version of the interval labeling algorithm of [1] tailored for *ST-graphs*. It can be shown that for a random string, the *ST-graph* can be preprocessed in $O(n \log n)$ expected time and space to answer reachability queries in $O(\log n)$ time. Furthermore, we present a series of strings for which their *ST-graphs* require $\Theta(n\sqrt{n})$ time and space of preprocessing when the algorithm is applied. Exhaustive computational calculations indicate that the series gives the worst case instances of the algorithm for the strings of length up to 33, strongly supporting that the worst case complexity of the preprocessing is $\Theta(n\sqrt{n})$ time and space. Assuming this is true, this would break the $O(n^2)$ barrier for total time and space used when conducting $O(n)$ queries.

Reachability on *ST-graphs* solves the problem of whether or not the string represented by the path from the root to the given query nodes are substrings of each other. The algorithm has possible applications in substring pattern set discovery, where the objective is to find best set of substrings that characterizes a given set (or sets) of strings: Consider two substring patterns such that one is a substring of the other. The set of strings which contain the former pattern as a substring is obviously a subset of the set of strings which have the latter pattern as a substring.

This property may allow us choose the best pattern set more efficiently, for example, by quickly detecting non-interesting pattern sets.

2. Preliminaries

Let Σ be a finite *alphabet*. An element of Σ^* is called a *string*. Strings x , y , and z are said to be a *prefix*, *substring*, and *suffix* of string $w = xyz$, respectively. The length of a string w is denoted by $|w|$. The empty string is denoted by ε , that is, $|\varepsilon| = 0$. Unless otherwise noted, we shall only consider strings of a fixed alphabet. Also, we assume that all strings end with a unique character $\$ \in \Sigma$ that does not occur anywhere else in the strings.

A *suffix tree* $\text{suftree}(s)$ for a given string s is a rooted tree whose edges are labeled with non-empty substrings of s , satisfying the following characteristics. For any node v in the suffix tree, let $\text{path}(v)$ denote the string spelled out by concatenating the edge labels on the path from the root to v . For each leaf node v , $\text{path}(v)$ is a distinct suffix of s , and for each suffix of s , there exists such a leaf v . Furthermore, each internal node has at least two children, and the first character of the labels on the edges to its children are distinct. The parent of node v is denoted by $\text{parent}_s(v)$, and the set of children of node v is denoted by $\text{children}_s(v)$. The *length* of node v is defined to be $|\text{path}(v)|$. The *depth* of node v with respect to the suffix tree is the number of edges on the path from the root to the node, and is always less than or equal to $|\text{path}(v)|$. The *height* of a suffix tree is the maximum depth of all nodes with respect to the suffix tree. Also, let $\text{subtree}_s(v)$ be the subtree of the suffix tree rooted at node v .

For a node v where $\text{path}(v) = \sigma x$ for some $\sigma \in \Sigma$ and $x \in \Sigma^*$, we denote the *suffix link* of v as $\text{parent}_l(v) = u$ where $\text{path}(u) = x$. It is easy to see that a unique $\text{parent}_l(v)$ exists for each node v in $\text{suftree}(s)$, except for the root node. Therefore, the suffix links also form a tree structure, which we denote by $\text{suflinktree}(s)$. Let $\text{children}_l(v) = \{u : \text{parent}_l(u) = v\}$. Note that the depth of node v with respect to the suffix link tree is always equivalent to $|\text{path}(v)|$. Also, let $\text{subtree}_l(v)$ be the subtree of the suffix link tree rooted at node v .

2.1. ST-graphs

Let V be the set of nodes of $\text{suftree}(s)$. Let us denote the set of (backward) edges of $\text{suftree}(s)$ by $E_s = \{(v, \text{parent}_s(v)) : v \in V\}$ and the set of edges of the $\text{suflinktree}(s)$ by $E_l = \{(v, \text{parent}_l(v)) : v \in V\}$. We define the *ST-graph* of a string s as the directed graph $G = (V, E)$ where $E = E_s \cup E_l$. It is well known that the suffix tree and its suffix links for a string of length n can be constructed and represented in $O(n)$ time and space [10, 6, 8, 4]. Fig. 1 shows an example of an *ST-graph* for the string $\text{ababbabbba\$}$. It is easy to see that the graph is not planar nor a partial lattice.

The problem we shall consider in this paper is as follows:

4 Yasuto Higa et al.

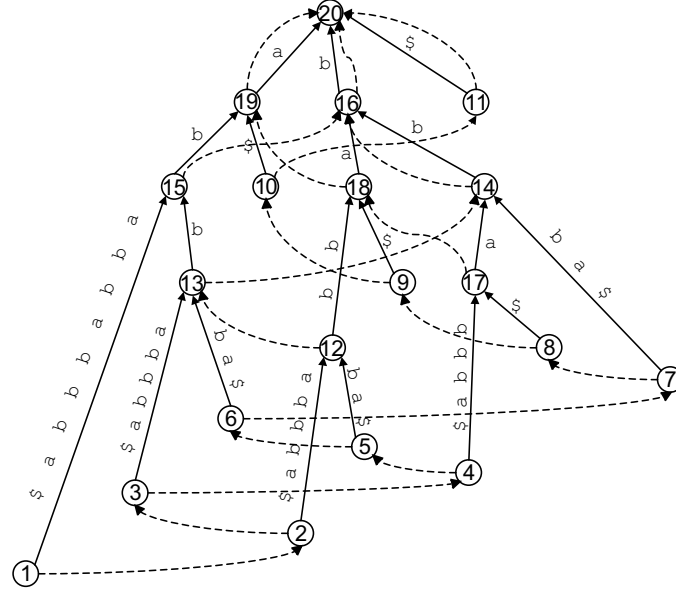


Fig. 1. A graph induced from the suffix tree of string **ababbabba\$**. Solid edges represent edges of the suffix tree. The dashed edges represent suffix links.

Table 1. Post-order interval and labels assigned to *ST*-graph of Fig. 1 by Algorithm 1.

node	interval	labels	node	interval	labels
1	[1,1]	[1,1]	11	[1,11]	[1,11]
2	[1,2]	[1,2]	12	[12,12]	[1,5],[12,12]
3	[1,3]	[1,3]	13	[12,13]	[1,6],[12,13]
4	[1,4]	[1,4]	14	[12,14]	[1,8],[12,14]
5	[1,5]	[1,5]	15	[15,15]	[1,6],[12,13],[15,15]
6	[1,6]	[1,6]	16	[12,16]	[1,9],[12,16],[17,18]
7	[1,7]	[1,7]	17	[17,17]	[1,8],[17,17]
8	[1,8]	[1,8]	18	[17,18]	[1,9],[12,12],[17,18]
9	[1,9]	[1,9]	19	[17,19]	[1,10],[12,13],[15,15],[17,19]
10	[1,10]	[1,10]	20	[1,20]	[1,20]

Problem 1 (*ST*-graph reachability query) Given the *ST*-graph $G = (V, E)$ of string s and an arbitrary pair of node $u, v \in V$, $rquery(u, v)$ returns **true** if there exists a path from node u to v in G , and **false** otherwise.

The query $rquery(u, v)$ is equivalent to the query of whether the string $path(v)$ is a substring of $path(u)$ or not.

Lemma 1. *Given an ST-graph $G = (V, E)$ of string s and nodes $u, v \in V$,*

$$rquery(u, v) = \text{true} \iff path(v) \text{ is a substring of } path(u)$$

Proof. (\Rightarrow) Suppose v is reachable from u . An edge $(p, q) \in E_s$ implies that $path(q)$ is a substring (prefix) of $path(p)$. An edge $(p, r) \in E_l$ implies that $path(r)$ is a substring (suffix) of $path(p)$. Since any path from u to v consists of edges in $E_s \cup E_l$, this implies that v is a substring of u .

(\Leftarrow) Suppose $path(v)$ is a substring of $path(u)$, i.e. there exists $x, z \in \Sigma^*$ such that $path(u) = xpath(v)z$. If $x = \varepsilon$, then $path(v)$ is a prefix of $path(u)$ which implies $u \in subtree_s(v)$, and that v is reachable from u using edges of E_s . For $x = x_1 \cdots x_k$ ($k \geq 1$), the nodes reachable from u using edges in E_l will have corresponding paths: $x_2 \cdots x_k path(v)z$, $x_3 \cdots x_k path(v)z$, \dots , $path(v)z$. Let v' be the node where $path(v') = path(v)z$. Then, since $path(v)$ is a prefix of $path(v')$, v is reachable from v' , and therefore reachable from u . \square

Corollary 2. *For any two distinct nodes $u, v \in V$ in an ST-graph (V, E) , if there exists a path from u to v , then there exists a path: $v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_t$ where $t > 0$, $u = v_0$, $v_t = v$, and for some $0 \leq i \leq t$, $(v_{j-1}, v_j) \in E_l$ for all $1 \leq j \leq i$ and $(v_{j-1}, v_j) \in E_s$ for all $i+1 \leq j \leq t$.*

Proof. Since there exists a path from u to v , $path(v)$ is a substring of $path(u)$ by Lemma 1. The Corollary follows from the argument for the (\Leftarrow) part of the proof of Lemma 1. \square

3. Interval Labeling Algorithm

In this paper, we analyze the complexity of the interval labeling algorithm for general DAGs presented in [1], when applied to ST-graphs. The original algorithm works as follows: First, create a spanning tree of the DAG by examining the nodes in topological order. The parent of each node is chosen so that the number of ancestors of each node is the largest. Then, intervals based on a postorder numbering of the spanning tree are assigned to each node. Further, the intervals of each node are propagated to all its ancestor nodes of the DAG, and as a result, each node will hold a set of interval labels. During the propagation, for a given set of interval labels at each node, *redundant* intervals which are subsumed by larger intervals in the same set are removed.

The interval labeling algorithm of [1] modified for ST-graphs is given in Algorithm 1, and Algorithm 2 shows how to answer $rquery(,)$ using the labels. Algorithm 1 has been simplified as follows: First, the topological ordering and spanning tree construction is not necessary. This is because each node v will only have at most two parents, one from the suffix tree, and the other from the suffix link tree. It is easy to verify that the path from the root to v using the edges of the suffix link tree is always at least as long as the path using the edges of the suffix tree.

Therefore, the suffix link tree already corresponds to the desired spanning tree, by which intervals are assigned to each node based on a postorder numbering. Second, although the original interval labeling algorithm requires labels to be propagated across all edges, this is not required for ST -graphs, and are propagated only across edges in E_s .

The correctness of the algorithm can be proved as follows. Suppose node v is reachable from node u . Then, node v must hold an interval label which subsumes the interval of u for Algorithm 2 to correctly answer $rquery(u, v)$. From lines 5-8 in Algorithm 1, node v holds all intervals of nodes in $subtree_s(v)$, which includes the interval for node v_i that can be reached by traversing suffix links starting from u as in Corollary 2. Since the interval of v_i subsumes the interval of u defined by the postorder numbering assigned in lines 1-4, v will contain an interval label which subsumes the interval of u .

4. Complexity

In this section, we will derive estimates on the complexity of Algorithm 1. In particular, we will consider the expected running time, as well as lower bounds for the worst case.

Lemma 3. *Assuming a constant size alphabet, Algorithm 1 runs in time linear in the total number of labels assigned to the nodes of the ST -graph.*

Proof. The maximum in-degree of each node is bounded by $O(|\Sigma|)$. This is because for the suffix tree, all labels on the edges must begin with a different character of the alphabet. For the suffix link tree, $|\{v : path(v) = \sigma path(u), \sigma \in \Sigma\}| \leq |\Sigma|$ for any node u . Therefore, assuming that $|\Sigma|$ is constant, merging the sorted labels (and removing subsumed intervals) from the in-coming nodes can be done in time linear in the total size of the in-coming labels. \square

From Lemma 3, we have only to count the number of labels that will be assigned to the ST -graph by Algorithm 1 in order to estimate the complexity of the algorithm.

4.1. Expected running time

A simple bound relating the height of the suffix tree and the total number of labels assigned to the ST -graph is shown in the following lemma.

Lemma 4. *For an ST -graph for a string of length n whose suffix tree has height h , the total number of interval labels assigned by Algorithm 1 is at most $O(nh)$.*

Proof. Notice that since the interval labels are only propagated through edges of the suffix tree, the maximum number of labels at a given node v is bounded by the number of nodes in $subtree_s(v)$. Therefore, at a given depth of the suffix tree, there can only be a maximum total of $O(n)$ labels, i.e., the total number of nodes

Algorithm 1: Assign labels to each node of the ST -graph.

Input: ST -graph $G(V, E)$ of string s
Output: Labeled ST -graph ($v.int$ and $v.labels$ for all $v \in V$)

```

1 foreach node  $v \in V$  in post-order of suflinktree(s) do
2    $v.int := [\min\{\text{post-order number of } subtree_l(v)\}, \text{post-order number of } v];$ 
3    $v.labels := \{v.int\};$ 
4 endfch
5 foreach node  $v \in V$  in post-order of suftree(s) do
6    $v.labels := \text{merge and sort } v.labels \text{ and } \{c.labels : c \in children_s(v)\};$ 
7   Remove  $[i, j] \in v.labels$  if  $\exists [i', j'] \in v.labels$  s.t.  $i' \leq i \leq j \leq j'$ ;
8 endfch

```

Algorithm 2: $rquery(u, v)$ on ST -graphs using labels assigned by Algorithm 1.

Input: Labeled ST -graph $G(V, E)$ of string s and nodes $u, v \in V$
Output: $rquery(u, v)$

```

1  $[i, j] = u.int;$ 
2 if  $\exists [i', j'] \in v.label$  such that  $i' \leq i \leq j \leq j'$  then return true;
3 return false;

```

in the suffix tree, since the subtrees of nodes of the same depth cannot intersect. This results in a maximum total of $O(nh)$ labels for all nodes. \square

Theorem 5. *The expected running time of Algorithm 1 for a random string of length n is $O(n \log n)$.*

Proof. It is known that the expected height of the suffix tree of a random string of length n is $O(\log n)$ [2]. The theorem follows from Lemma 3 and Lemma 4. \square

4.2. Worst Case Lower Bounds

In this subsection, we will give a lower bound for the worst case complexity of Algorithm 1. We will present a series of strings of length n whose ST -graphs will have $\Theta(n\sqrt{n})$ labels assigned by the algorithm. Prior to this, we show related properties of suffix trees and suffix link trees.

4.2.1. Properties of ST -graphs.

For $|\Sigma| = 2$, there can only be one string of a given length n , and the structure of the ST -graph is determined uniquely (recall that all strings terminate with a uniquely occurring character $\$ \in \Sigma$). It is not difficult to show that the total number of labels

is $3(n-1) = O(n)$ in such case. In what follows we therefore consider the case for $|\Sigma| \geq 3$.

Lemma 6. *The number of interval labels assigned to each node of the ST-graph for any string s is bounded by the number of leaves in $\text{suflinktree}(s)$. More generally, the exact number of labels assigned to node v is $\min_{W \subseteq \text{subtree}_s(v)} \{|W| : \text{subtree}_s(v) \subseteq \cup_{w \in W} \text{subtree}_l(w)\}$.*

Proof. Let ℓ be the number of leaves in $\text{suflinktree}(s)$. In the post-order traversal on $\text{suflinktree}(s)$ (lines 1-4 of Algorithm 1), each node receives exactly one interval label, and there are exactly ℓ different values for the first element of the intervals. In post-order traversal on $\text{suftree}(s)$ (lines 5-8 of Algorithm 1), we remove subsumed intervals and therefore each node gets at most ℓ interval labels. The exact number of labels follows from a similar argument. \square

Lemma 7. *If and only if $|\text{children}_s(v)| = 2$ for any internal node $v \neq \text{root}$ of $\text{suftree}(s)$, the number of internal nodes (excluding the root) is the maximum, which is $n-3$.*

Proof. Because of $\$$ there are always n leaves in $\text{suftree}(s)$. Since there is always an edge labeled with $\$$ from the root, $|\text{children}_s(\text{root})| = |\Sigma|$. \square

Lemma 8. *Assume $\text{suftree}(s)$ satisfies the condition in Lemma 7. For the three following groups of the internal nodes of $\text{suftree}(s)$,*

- n_1 : internal nodes with two leaf-children;
- n_2 : internal nodes with one leaf-child and one internal-child;
- n_3 : internal nodes with two internal-children;

where $n_1 + n_2 + n_3 = n-3$, we have

$$\begin{aligned} n_2 &= n - 2n_1 - 1, \text{ and} \\ n_3 &= n_1 - 2. \end{aligned}$$

Proof. Because of $\$$, $\text{suftree}(s)$ always has n leaves. Since there are n_1 internal nodes with two leaf-children, and since the leaf corresponding to suffix $\$$ is a child of the root, we have $n_2 = n - 2n_1 - 1$. Finally $n_3 = n - 3 - n_1 - n_2 = n_1 - 2$. \square

Lemma 9. *Assume $\text{suftree}(s)$ satisfies the condition in Lemma 7. For any node v , if $|\text{children}_l(v)| \geq 2$, then v is a group n_3 node in $\text{suftree}(s)$.*

Proof. Let $x = \text{path}(v)$, and let $u, w \in \text{children}_l(v)$. Since u and w are nodes in the suffix tree, there exists at least four possible suffixes whose paths must pass node v . We denote these paths as $x\sigma_1y_1, x\sigma_2y_2, x\sigma_3y_3, x\sigma_4y_4$ where $\sigma_i \in \Sigma$ and $y_i \in \Sigma^*$ for $1 \leq i \leq 4$. It must be that $\sigma_{i_1} = \sigma_{i_2}$ and $\sigma_{i_3} = \sigma_{i_4}$ for some $\{i_1, i_2, i_3, i_4\} =$

$\{1, 2, 3, 4\}$. Since all four paths must be distinct, it follows that there must exist distinct child nodes of v where $x\sigma_{i_1}y_{i_1}$ and $x\sigma_{i_2}y_{i_2}$ diverge, and $x\sigma_{i_3}y_{i_3}$ and $x\sigma_{i_4}y_{i_4}$ diverge respectively. \square

Lemma 10. *Assume $\text{suftree}(s)$ satisfies the condition in Lemma 7. The number ℓ of leaves in $\text{suflinktree}(s)$ is at most $n_1 + 1$.*

Proof. Let ℓ be the number of leaves in $\text{suflinktree}(s)$. Since all leaves of $\text{suftree}(s)$ are included in one path of suffix links and this path leads to the root, the maximum number of internal nodes v of $\text{suflinktree}(s)$, such that $|\text{children}_l(v)| \geq 2$, is $\ell - 3$ excepting the root. From Lemma 9, this leads to $\ell - 3$ internal nodes with two internal-children in $\text{suftree}(s)$. Therefore, $\ell \leq n_3 + 3 = n_1 + 1$ by Lemma 8. \square

Lemma 11. *Assume $\text{suftree}(s)$ satisfies the condition in Lemma 7. If $\ell = n_3 + 3$, then the longest internal node in group n_3 has at most $\ell - 1$ labels.*

Proof. It follows from Lemmas 8, 9, and 10. \square

4.2.2. Lower Bound.

Consider the following series of strings of length $n = \frac{3}{2}i(i+3)+5$ where $i = 1, 2, 3, \dots$:

$$X_i = \text{ab}^i \text{ab}^{i+1} \text{ab}^i \text{ab}^1 \text{ab}^i \text{ab}^{i-1} \text{ab}^i \text{ab}^2 \text{ab}^i \text{ab}^{i-2} \text{a} \dots \text{ab}^k \text{ab}^i \text{ab}^{i-k} \text{a} \dots \text{ab}^i \text{ab}^{\lceil i/2 \rceil} \text{ab}^i \text{ab}^i \text{a}\$$$

In what follows, we analyze X_i in terms of internal nodes of $\text{suftree}(X_i)$, shown in Fig. 2. We consider the structure of $\text{suftree}(X_i)$ in Lemmas 13, 14, 15, 17, and the structure of $\text{suflinktree}(X_i)$ in Lemmas 18, 19, 20, 21. An important point in the lemmas is that if a substring p appears only once in the string, then there is no explicit internal node v such that $\text{path}(v) = p$, and the substring will correspond to a position on an edge leading to a leaf node (or the leaf node itself). Also, if $p\sigma_1$ and $p\sigma_2$ are both substrings of the string for distinct $\sigma_1, \sigma_2 \in \Sigma$, then there exists a node v such that $\text{path}(v) = p$.

Lemma 12. *$\text{suftree}(X_i)$ satisfies the condition in Lemma 7.*

Proof. Any occurrence of a in X_i is followed by either b or $\text{\$}$, and any occurrence of b in X_i is followed by either a or b . Moreover, $\text{\$}$ appears only at the end of X_i . Thus, for any internal node v of $\text{suftree}(X_i)$, we have $|\text{children}_s(v)| = 2$. \square

Lemma 13. *For any k ($0 \leq k \leq i$), there exists an internal node corresponding to $\text{b}^{i-k} \text{ab}^i \text{ab}^k$, and this node belongs to group n_1 of $\text{suftree}(X_i)$.*

Proof. We have three cases to consider:

- When $k = 0$.

Since there are two strings $\text{b}^i \text{ab}^i \text{ab}$ and $\text{b}^i \text{ab}^i \text{a}\$$ in X_i , there is an internal

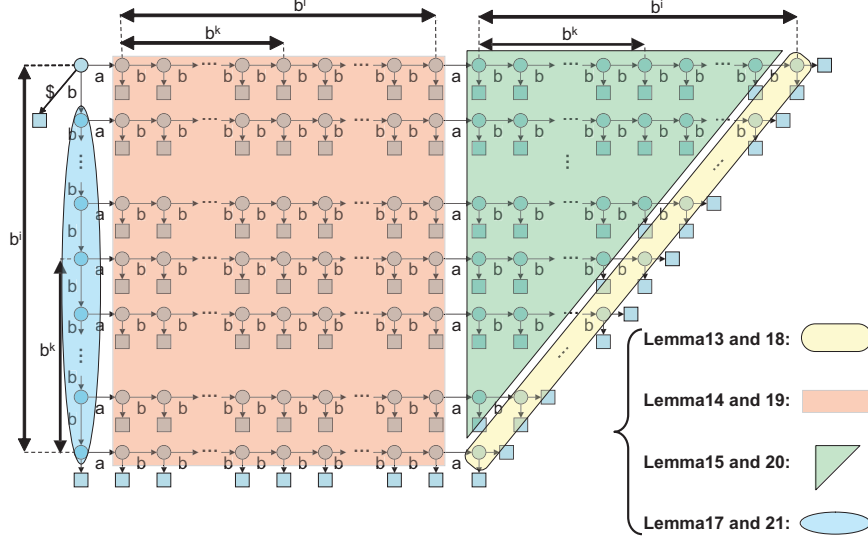


Fig. 2. Illustration for $\text{suftree}(X_i)$. The four groups of the internal nodes are dealt in Lemmas 13, 14, 15, and 17, respectively.

node for $b^i a b^i a$. Since there is no other occurrence of $b^i a b^i a$, the two children of this node are both a leaf node and thus it belongs to group n_1 .

- When $0 < k < i$.

Consider a substring Y_i of X_i such that

$$Y_i = ab^{k-1} \underline{ab^i ab^{i-k+1} ab^i ab^k ab^i ab^{i-k} ab^i ab^{k+1} ab^i a}.$$

Then, $b^{i-k} ab^i ab^k$ appears twice in this substring, as underlined.

Now we show that each of $b^{i-k} ab^i ab^k a$ and $b^{i-k} ab^i ab^k b$ appears only once in X_i . For $b^{i-k} ab^i ab^k a$, it is clear because $ab^k a$ appears exactly once in X_i . String $b^{i-k} ab^i ab^k b = b^{i-k} ab^i ab^{k+1}$ appears in Y_i (the second underlined part). This is the only occurrence of $b^{i-k} ab^i ab^{k+1}$ in X_i , because the prefix $b^{i-k} a$ would appear in substrings $ab^{i-k+x} ab^i ab^{k-x} a$ with $x \geq 0$, but then the suffix $b^i ab^{k+1}$ cannot match.

- When $k = i$.

By similar discussions to the case that $k = 0$. □

Lemma 14. For any x ($0 \leq x \leq i$) and y ($0 \leq y \leq i$), there exists an internal node corresponding to $b^x a b^y$, and this node belongs to group n_2 .

Proof. We have three cases to consider:

- When $y = 0$.

Since there are more than one occurrences of $b^i a b$, there are more than one

occurrences of $\mathbf{b}^x \mathbf{ab}$. In addition, since there is exactly one occurrence of $\mathbf{b}^i \mathbf{a} \$$, there is exactly one occurrence of $\mathbf{b}^x \mathbf{a} \$$.

- When $0 < y < i$.

Since there are more than one occurrences of $\mathbf{b}^x \mathbf{ab}^i$, there are more than one occurrences of $\mathbf{b}^x \mathbf{ab}^{y+1}$. In addition, since there is exactly one occurrence of $\mathbf{b}^i \mathbf{ab}^y \mathbf{a}$ for each $0 < y < i$, $\mathbf{b}^x \mathbf{ab}^y \mathbf{a}$ appears exactly once.

- When $y = i$.

There is exactly one occurrence of $\mathbf{b}^x \mathbf{ab}^i \mathbf{a}$ for each $0 \leq x \leq i$. Since there is exactly one occurrence of $\mathbf{b}^i \mathbf{ab}^{i+1}$, $\mathbf{b}^x \mathbf{ab}^{i+1}$ appears exactly once for each $0 \leq x \leq i$. \square

Lemma 15. *For any y ($0 \leq y < k$), where $1 \leq k \leq i$, there exists an internal node corresponding to $\mathbf{b}^{i-k} \mathbf{ab}^i \mathbf{ab}^y$, and this node belongs to group n_2 .*

Proof. By similar arguments to Lemmas 13 and 14. \square

The next corollary follows Lemmas 14 and 15.

Corollary 16. *For any k ($0 \leq k \leq i$), let $\text{path}(u) = \mathbf{b}^{i-k} \mathbf{a}$ and $\text{path}(v) = \mathbf{b}^{i-k} \mathbf{ab}^i \mathbf{ab}^k$. Then, the path from u to v contains $i + k + 2$ internal nodes (including u and v).*

Lemma 17. *For any x ($0 < x \leq i$), there is an internal node corresponding to \mathbf{b}^x . The node for \mathbf{b}^i belongs to group n_2 , and the other nodes for \mathbf{b}^x with $0 < x < i$ belong to group n_3 .*

Proof. Since $\mathbf{b}^i \mathbf{a}$ and \mathbf{b}^{i+1} appear in \mathbf{X}_i , there is an internal node for \mathbf{b}^x for any $0 < x \leq i$. Since $\mathbf{b}^i \mathbf{a}$ appears more than once and \mathbf{b}^{i+1} appears exactly once, the node for \mathbf{b}^i belongs to group n_2 . All the nodes for \mathbf{b}^x with $0 < x < i$ belong to group n_3 , since both $\mathbf{b}^x \mathbf{a}$ and \mathbf{b}^{x+1} appear more than once. \square

From here on, we consider the suffix links of $\text{suftree}(\mathbf{X}_i)$.

Lemma 18. *For any k ($0 \leq k \leq i$), let v be the internal node such that $\text{path}(v) = \mathbf{b}^{i-k} \mathbf{ab}^i \mathbf{ab}^k$. Then we have $|\text{children}_l(v)| = 0$.*

Proof. For contrary, assume that $|\text{children}_l(v)| \geq 1$. Then there must exist a node corresponding to either $\mathbf{ab}^{i-k} \mathbf{ab}^i \mathbf{ab}^k$ or $\mathbf{b}^{i-k+1} \mathbf{ab}^i \mathbf{ab}^k$. However, we show that these strings can appear at most once in \mathbf{X}_i .

First, we consider $\mathbf{ab}^{i-k} \mathbf{ab}^i \mathbf{ab}^k$:

- When $k = 0$. $\mathbf{ab}^i \mathbf{ab}^i \mathbf{a}$ appears only once (at the end of \mathbf{X}_i).
- When $0 < k < i$. Prefix $\mathbf{ab}^{i-k} \mathbf{a}$ appears only once in \mathbf{X}_i .
- When $k = i$. Prefix \mathbf{aa} never appears in \mathbf{X}_i .

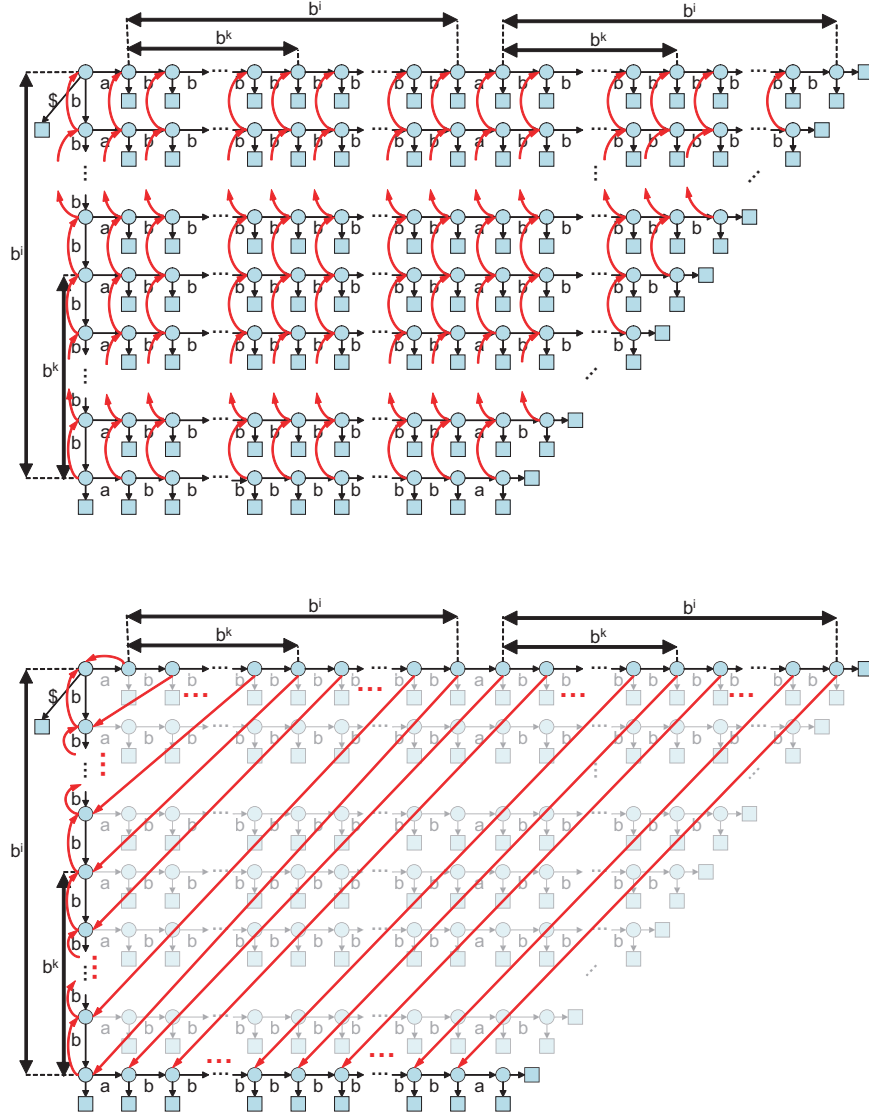


Fig. 3. Illustration for the suffix links of $\text{suftree}(X_i)$. For the sake of visibility, the suffix links of X_i are shown in two rounds. Moreover, for simplicity, the suffix links for the leaves are omitted here.

Now let us consider $b^{i-k+1}ab^iab^k$. By Lemma 13, there is an internal node corresponding to $b^{i-k+1}ab^iab^{k-1}$ and this node belongs to group n_1 . This implies that $b^{i-k+1}ab^iab^k$ can appear at most once in X_i . \square

The above lemma implies that the internal nodes v such that $\text{path}(v) = b^{i-k}ab^iab^k$ are the leaves of $\text{sufflinktree}(X_i)$. See also the upper diagram in Fig. 3.

Lemma 19. *For any internal node $v_{x,y}$ such that $\text{path}(v_{x,y}) = \mathbf{b}^x \mathbf{a} \mathbf{b}^y$ ($0 \leq x \leq i$ and $0 \leq y \leq i$), we have $|\text{children}_l(v_{x,y})| = 1$.*

Proof. We have three cases to consider:

- When $x = 0$.
There is no occurrence of $\mathbf{a} \mathbf{a} \mathbf{b}^y$, and we have two distinct occurrences of $\mathbf{b}^i \mathbf{a} \mathbf{b}^i$ in \mathbf{X}_i . Thus we have $|\text{children}_l(v_{0,y})| = 1$.
- When $0 < x < i$.
 $\mathbf{a} \mathbf{b}^x \mathbf{a} \mathbf{b}^y$ appears only once, since there is only one occurrence of $\mathbf{a} \mathbf{b}^x \mathbf{a} \mathbf{b}^i$. Because $\mathbf{b}^i \mathbf{a} \mathbf{b}^y$ appears more than once, there are more than one occurrence of $\mathbf{b}^x \mathbf{a} \mathbf{b}^y$. Thus we have $|\text{children}_l(v_{x,y})| = 1$.
- When $x = i$.
We have at least two occurrences of $\mathbf{a} \mathbf{b}^i \mathbf{a} \mathbf{b}^y$. Since \mathbf{b}^{i+1} appears only once, there is only one occurrence of $\mathbf{b}^{i+1} \mathbf{a} \mathbf{b}^y$. Thus we have $|\text{children}_l(v_{i,y})| = 1$.

Lemma 20. *For any internal node v_z such that $\text{path}(v_z) = \mathbf{b}^{i-k} \mathbf{a} \mathbf{b}^i \mathbf{a} \mathbf{b}^z$ ($0 \leq z < k$), where $1 \leq k \leq i$, we have $|\text{children}_l(v_z)| = 1$.*

Proof. By similar arguments to Lemma 19. □

Lemma 21. *For any internal node v_x such that $\text{path}(v_x) = \mathbf{b}^x$ ($0 < x < i$), we have $|\text{children}_l(v_x)| = 2$. In addition, for the node v_i such that $\text{path}(v_i) = \mathbf{b}^i$, we have $|\text{children}_l(v_i)| = 1$.*

Proof. Due to Lemma 14, there exist internal nodes u_x such that $\text{path}(u_x) = \mathbf{a} \mathbf{b}^x$ for each $0 < x \leq i$, and therefore have $(u_x, v_x) \in E_l$. Due to Lemma 17, we have $(v_{x+1}, v_x) \in E_l$ for $0 < x < i$. For $x = i$, however, we have $(v_{i+1}, v_i) \notin E_l$ because \mathbf{b}^{i+1} does not correspond to a node (since it occurs only once in the string). □

From the above lemmas, only the nodes corresponding to \mathbf{b}^x ($0 < x < i$) are of in-degree two in $\text{suflinktree}(\mathbf{X}_i)$. Plus, only the root node is of in-degree three in $\text{suflinktree}(\mathbf{X}_i)$. See also Fig. 3 for these observations.

The number of nodes covered in these lemmas are as follows:

Lemma	11 or 15	12 or 16	13 or 17	14 or 18
The number of nodes	$i + 1$	$(i + 1)^2$	$i(i + 1)/2$	i

the sum of these nodes is $\frac{3}{2}i^2 + \frac{9}{2}i + 2 = n - 3$, which indicates that we have discussed all nodes of the ST -graph for the string \mathbf{X}_i . Now we are finally ready to show the lower bound on the number of interval labels.

Theorem 22. *The number of labels assigned to a suffix tree by Algorithm 1 is $\Omega(n\sqrt{n})$.*

Proof. From Lemmas 6, 11, 12, 13, 14, 15, 17, 18, 19, 20, 21, Corollary 16, we have $\ell = i + 2$ and achieve the following equations on the total number of interval labels that are assigned to $\text{suftree}(\mathbf{X}_i)$:

$$\begin{aligned}
 & (2 + 3 + \cdots + \ell - 1) \times (i + 1) + \ell \times \left\{ 1 + \sum_{k=0}^i \{(i + k + 2) - (\ell - 2)\} \right\} \\
 & + \{(\ell - 1) + (\ell - 2) + \cdots + 3\} + 1 \times (n + 1) \\
 = & \{2 + 3 + \cdots + (i + 1)\} \times (i + 1) + (i + 2) \times \left\{ 1 + \sum_{k=0}^i (k + 2) \right\} \\
 & + \{(i + 1) + i + \cdots + 3\} + \frac{3i(i + 3)}{2} + 6 \\
 = & \frac{i(i + 3)(i + 1)}{2} + \frac{(i + 2)\{2 + i(i + 1) + 4(i + 1)\}}{2} + \frac{(i - 1)(i + 4)}{2} + \frac{3i(i + 3)}{2} + 6 \\
 = & \frac{2i^3 + 15i^2 + 31i + 20}{2}. \tag{1}
 \end{aligned}$$

Since $n = \frac{3}{2}i(i + 3) + 5$ and $i \geq 1$, we have

$$i = \frac{\sqrt{24n - 39}}{6} - \frac{3}{2}.$$

By substituting this for the i 's in Equation (1), the total number of interval labels assigned to $\text{suftree}(\mathbf{X}_i)$ is shown to be $\Theta(n\sqrt{n})$. \square

The worst case upper bound for Algorithm 2 is $O(\log n)$: From the argument in Lemma 3, the labels at each node can be stored as sorted arrays. Also, the maximum number of labels at each node is bounded by $O(n)$ (Lemma 6). Therefore, line 2 in Algorithm 2 can be run in $O(\log n)$ time using a standard binary search on the label array.

5. Computational Experiments

We exhaustively enumerated all strings of length $n \leq 33$ consisting of $\{\mathbf{a}, \mathbf{b}\}$ and ending with $\mathbf{\$}$, and applied Algorithm 1 to each string. For each n , the number of labels in the worst case was recorded. The results are shown in Table 2.

We note that for $n \leq 7$, the worst case corresponds to the string $\mathbf{a}^{n-1}\mathbf{\$}$, and therefore the total number of labels is $3(n - 1)$. For $n = 11, 20, 32$, the worst case instances contain \mathbf{X}_1 , \mathbf{X}_2 , \mathbf{X}_3 , and the total number of labels is as calculated in Theorem 22. Generally for $7 \leq n \leq 33$, we found that the total number of labels in the worst case can be written exactly with the following formula:

$$\max\{f(\lfloor \sqrt{(2n - 3)/3} \rfloor), f(\lceil \sqrt{(2n - 3)/3} \rceil)\}$$

where $f(k) = (k + 2)n - k(k^2 + 3)/2 - 3$. We have also confirmed for smaller n and with larger sized alphabets, the worst case instances will only contain one type of character (excluding $\mathbf{\$}$) for $n \leq 7$, and two types of characters (excluding $\mathbf{\$}$) for

Table 2. The maximum number of labels that is assigned by Algorithm 1 to any string of length n .

n	2	3	4	5	6	7	8	9	10	11	12
#labels	3	6	9	12	15	18	22	26	30	34	39
n	13	14	15	16	17	18	19	20	21	22	23
#labels	44	49	54	59	64	69	74	79	85	91	97
n	24	25	26	27	28	29	30	31	32	33	
#labels	103	109	115	121	127	133	139	145	151	158	

$n \geq 7$, therefore corresponding to the instances given above. (At $n = 7$, both types had equal worst case label size of 18.) This seems natural since Lemma 7 indicates that the more types of characters used, the less number of nodes there are in the ST -graph.

Although we have not been able to give a rigorous proof for an upperbound of $O(n\sqrt{n})$, the above results strongly suggest this bound.

6. Discussion

We presented an algorithm that can process an ST -graph for a string of length n , so that reachability queries between arbitrary pairs of nodes can be answered in $O(\log n)$ time. The expected time and space complexity of the preprocessing algorithm for a random string is $O(n \log n)$. We also presented a series of strings for which the algorithm requires $\Theta(n\sqrt{n})$ time and space for preprocessing. Exhaustive computational search for $n \leq 33$ showed that the strings of the series also achieve the worst case instances of the algorithm. Although we have not been able to give a direct proof, this provides strong evidence that the worst case time complexity of the algorithm is also $\Theta(n\sqrt{n})$.

Since a suffix tree can have a height of $O(n)$, a naïve consideration of Lemma 4 only gives an $O(n^2)$ bound for the number of labels for an ST -graph of a suffix tree of height $O(n)$, rather than $O(n\sqrt{n})$. There seems to be a delicate tradeoff between deep paths in the suffix tree and deep paths in the suffix link tree. For example, the suffix tree for string $a^n\$$ will have a path of depth $O(n)$. However, the number of total labels in this case is also limited to $O(n)$, since there are only two leaves in the suffix link tree, which bounds the number of labels for each node to 2 (Lemma 6). There also exists strings $a^k \underbrace{b b c c d d e e f f \dots}_k \$$ ($n = 3k + 1$), where their suffix trees have a path of depth $O(n)$, and their suffix link trees have $O(n)$ leaves. However, the total number of labels in this case is also $O(n)$, since all but two of the leaves in the suffix link tree are very shallow.

6.1. Open Problems

There are several open problems that are of interest concerning reachability queries on *ST*-graphs.

- (1) Whether or not there exists an algorithm which can do better: $O(n\sqrt{n})$ preprocessing and $O(1)$ query, or ultimately $O(n)$ preprocessing and $O(1)$ query.
- (2) Whether or not we can simulate reachability queries on suffix *trie* graphs. All nodes in the suffix tree correspond to a substring of the original string. However, there can exist substrings in the string without a corresponding explicit node. An implicit node of a suffix tree is a position in the suffix tree which ends somewhere in the middle of an edge. So far we have considered reachability queries between explicit nodes of the *ST*-graph. Although reachability between implicit nodes of the *ST*-graph is straightforward with respect to the *ST*-graph, the result does not correspond to the substring relation between implicit nodes, as it does for explicit nodes shown in Lemma 1.

References

- [1] R. AGRAWAL, A. BORGIDA, AND H. V. JAGADISH: *Efficient management of transitive relationships in large data and knowledge bases*, in Proceedings of the 1989 ACM SIGMOD International Conference on Management of Data, 1989, pp. 253–262.
- [2] A. APOSTOLICO AND W. SZPANKOWSKI: *Self-alignments in words and their applications*. Journal of Algorithms, 13(3) 1992, pp. 446–467.
- [3] E. COHEN, E. HALPERIN, H. KAPLAN, AND U. ZWICK: *Reachability and distance queries via 2-hop labels*, in Proc. 13th Annual ACM-SIAM Symposium on Discrete Algorithms, 2002, pp. 937–946.
- [4] D. GUSFIELD: *Algorithms on Strings, Trees, and Sequences*, Cambridge University Press, 1997.
- [5] T. KAMEDA: *On the vector representation of the reachability in planar directed graphs*. Information Processing Letters, 3(3) 1975, pp. 75–77.
- [6] E. M. MCCREIGHT: *A space-economical suffix tree construction algorithm*. J. ACM, 23(2) 1976, pp. 262–272.
- [7] M. TALAMO AND P. VOCCA: *A data structure for lattice representation*. Theoretical Computer Science, 175(2) 1997, pp. 373–392.
- [8] E. UKKONEN: *On-line construction of suffix trees*. Algorithmica, 14(3) 1995, pp. 249–260.
- [9] H. WANG, H. HE, J. YANG, P. S. YU, AND J. X. YU: *Dual labeling: answering graph reachability queries in constant time*, in Proc. 22nd International Conference on Data Engineering, 2006, to appear.
- [10] P. WEINER: *Linear pattern matching algorithms*, in Proc. 14th IEEE Annual Symp. on Switching and Automata Theory, 1973, pp. 1–11.