



ELSEVIER

Contents lists available at ScienceDirect

## Theoretical Computer Science

[www.elsevier.com/locate/tcs](http://www.elsevier.com/locate/tcs)
Ranked document selection <sup>☆</sup>J. Ian Munro <sup>a</sup>, Gonzalo Navarro <sup>b</sup>, Rahul Shah <sup>c</sup>, Sharma V. Thankachan <sup>d,\*</sup><sup>a</sup> Cheriton School of Computer Science, University of Waterloo, Canada<sup>b</sup> Center for Biotechnology and Bioengineering (CeBiB) & Millennium Institute for Foundational Research on Data (IMFD), Dept. of Computer Science, University of Chile, Chile<sup>c</sup> School of Electrical Engineering and Computer Science, Louisiana State University, USA<sup>d</sup> Department of Computer Science, University of Central Florida, USA

## ARTICLE INFO

## Article history:

Received 8 January 2019

Received in revised form 25 September 2019

Accepted 5 October 2019

Available online 9 October 2019

## Keywords:

Document Indexing

Top-k Document Retrieval

Succinct Data Structures

## ABSTRACT

Let  $\mathcal{D}$  be a collection of string documents of  $n$  characters in total. The *top-k document retrieval problem* is to preprocess  $\mathcal{D}$  into a data structure that, given a query  $(P, k)$ , can return the  $k$  documents of  $\mathcal{D}$  most relevant to pattern  $P$ . The relevance of a document  $d$  for a pattern  $P$  is given by a predefined ranking function  $w(P, d)$ . Linear space and optimal query time solutions already exist for this problem. In this paper we consider a novel problem, *document selection*, in which a query  $(P, k)$  aims to report the  $k$ th document most relevant to  $P$  (instead of reporting all top- $k$  documents). We present a data structure using  $O(n \log^\epsilon n)$  space, for any constant  $\epsilon > 0$ , answering selection queries in time  $O(\log k / \log \log n)$ , and a linear-space data structure answering queries in time  $O(\log k)$ , given the locus node of  $P$  in a (generalized) suffix tree of  $\mathcal{D}$ . We also prove that it is unlikely that a succinct-space solution for this problem exists with poly-logarithmic query time, and that  $O(\log k / \log \log n)$  is indeed optimal within  $O(n \text{ polylog } n)$  space for most text families. Finally, we present some additional space-time trade-offs exploring the extremes of those lower bounds.

© 2019 Elsevier B.V. All rights reserved.

## 1. Introduction

Document retrieval is a special branch of pattern matching related to information retrieval and web searching. In this problem, the data consists of a collection of text *documents*, and the queries refer to documents rather than text positions [23]. In this paper we focus on arguably the most important of those problems, called *top-k document retrieval*: Given  $\mathcal{D} = \{d_1, d_2, d_3, \dots, d_D\}$ , of total length  $n = \sum_{i=1}^D |d_i|$ , preprocess it into a data structure that, given a pattern  $P$  and a threshold  $k$ , retrieves the  $k$  documents from  $\mathcal{D}$  that are most *relevant* to  $P$ , in decreasing order of relevance. The relevance of a document  $d$  with respect to  $P$  is captured using any function  $w(P, d)$  of the starting positions of the occurrences of  $P$  in  $d$ . A popular example of relevance is the *term frequency* metric, that is, the number of occurrences of  $P$  in  $d$ . This is a well studied problem, and an existing linear space data structure can answer queries in optimal time  $O(k)$  [32], once the locus

<sup>☆</sup> Funded in part by NSERC of Canada and the Canada Research Chairs program; Fondecyt Grant 1-170048, Chile; Basal Funds FB0001, CONICYT, Chile; Millennium Institute for Foundational Research on Data (IMFD), Chile; and NSF Grants CCF-1017623, CCF-1218904. An early partial version of this paper appeared in *Proc. SWAT 2014* [20].

\* Corresponding author.

E-mail addresses: [imunro@uwaterloo.ca](mailto:imunro@uwaterloo.ca) (J. Munro), [gnavarro@dcc.uchile.cl](mailto:gnavarro@dcc.uchile.cl) (G. Navarro), [rahul@csc.lsu.edu](mailto:rahul@csc.lsu.edu) (R. Shah), [Sharma.Thankachan@ucf.edu](mailto:Sharma.Thankachan@ucf.edu) (S.V. Thankachan).

node of  $P$  in a generalized suffix tree of  $\mathcal{D}$  is found. Note that the locus of  $P$  can be computed in optimal  $O(|P|/\log_\sigma n)$  time in the RAM model [25], where  $\sigma$  is the alphabet size.

*Our Contribution.* In this paper we study a new related problem called *document selection*, where we must return the  $k$ th document of  $\mathcal{D}$  most relevant to  $P$ , that is, the  $k$ th element returned by a top- $k$  query (breaking ties arbitrarily). We present three main results, depending on the amount of space used: (1) We give a data structure that uses  $O(n \log^\epsilon n)$  space, for any constant  $\epsilon > 0$ , and answers queries in time  $O(\log k / \log \log n)$ , given the locus of  $P$ , (2) We give a linear-space data structure that answers queries in  $O(\log k)$  time, given the locus of  $P$ , (3) We prove that it is highly unlikely that the problem can be solved in less than linear space within poly-logarithmic time, via a reduction from the *position restricted substring searching* problem [17,11]. We also prove that the query time  $O(\log k / \log \log n)$  is indeed optimal within  $O(|\text{GST}| \text{polylog } n)$  space, where GST is the generalized suffix tree of  $\mathcal{D}$ , via a reduction from the *prefix selection* problem [15]. Lastly, we present some additional space-time trade-offs that explore the limits of these lower bounds.

*Advanced Queries.* The document selection is useful for various advanced queries. When a user browses ranked results of a (top- $k_1$ ) query and asks for the next set of (top- $k_2$ ) results with  $k_1 < k_2$ , we need to report the top- $k_2$  documents that are not top- $k_1$ . Instead of computing a top- $k_2$  query in time  $O(|P|/\log_\sigma n + k_2)$ , which may be non-optimal if  $\delta = k_2 - k_1 = o(k_2)$ , our results allow answering this query in time  $O(|P|/\log_\sigma n + \delta \log k_2)$  via  $\delta$  selection queries (note that the cost  $|P|/\log_\sigma n$  is paid only once). Another possible query is to count the number  $K$  of documents  $d$  with  $w(P, d) \geq \tau$ , given  $P$  and  $\tau$ . This can be answered via doubling search using document selection queries, in time  $O(|P|/\log_\sigma n + \log^2 K)$ , assuming  $w(P, d)$  can be computed in constant time given the locus of  $P$ . Similarly, we can count or list the documents  $d$  with  $w(P, d) \in [\tau_1, \tau_2]$ . Such queries are important in bioinformatics, for example for motif mining or for avoiding sequences where  $P$  is “over-expressed”, and for data mining in general, for example to estimate the distribution of relevance scores of certain patterns.

*Related Work.* The notion of relevance-based string retrieval was introduced by Muthukrishnan [22], who proposed and solved various problems but not top- $k$  document retrieval. The first data structure for this problem, under the term frequency measure and using  $O(n \log n)$  words of space, was given by Hon et al. [10]. Later, Hon et al. [12,13] introduced a linear space structure ( $O(n)$  words), that works for general weight functions as described earlier, with query time  $O(p + k \log k)$ . This was improved to  $O(p + k)$  [24,25], and finally to the optimal  $O(k)$  [32], all using linear space. Those times are in addition to the time for finding the locus node of  $P$ ,  $\text{locus}(P)$ , in the generalized suffix tree of  $\mathcal{D}$ , GST. The problem has also been studied in scenarios where less than linear space (i.e.,  $o(n \log n)$  bits) can be used. For example, it is possible to solve the problem efficiently using  $n \log \sigma + o(n \log \sigma)$  bits [27,33], where  $\sigma$  is the alphabet size of the text (thus  $n \log \sigma$  bits are needed to represent the text itself). The results are mostly tailored to the term frequency measure of relevance, and achieve times of the form  $O(k \text{polylog } n)$ . See [9,12,14,23,28] for more details. Later, Munro et al. [19,21] showed that an  $n \log \sigma + o(n \log \sigma)$  bits space index can be designed for term-proximity based relevance measures as well. Also see [1,12] for its extensions to the case where the query consists of two or more patterns.

## 2. The top- $k$ framework

This section briefly describes the linear-space framework of Hon et al. [12] for top- $k$  queries. The generalized suffix tree (GST) of a document collection  $\mathcal{D} = \{d_1, d_2, d_3, \dots, d_D\}$  is the combined compact trie of all the non-empty suffixes of all the documents [34]. The total number of leaves in GST is the same as the total length  $n$  of all the documents. For each node  $j$  in GST,  $\text{prefix}(j)$  is the string obtained by concatenating the edge labels on the path from the root to node  $j$ . The highest node  $v$  satisfying that  $P$  is a prefix of  $\text{prefix}(v)$  is called the *locus* of  $P$  and denoted  $\text{locus}(P) = v$ .

Let  $\ell_i$  represent the  $i$ th leftmost leaf node in GST. We say that a node is *marked* with a document  $d$  if it is either a leaf node whose corresponding suffix belongs to  $d$ , or it is the lowest common ancestor (LCA) of two such leaves. This implies that the number of nodes marked with document  $d$  is exactly equal to the number of nodes in the suffix tree of  $d$  (at most  $2|d|$ ). A node can be marked with multiple documents. For each node  $j$  and each of its marking documents  $d$ , define a *link* to be a quadruple  $\langle \text{doc} = d, \text{weight}, (\text{origin} = j, \text{target}) \rangle$ , where  $\text{weight} = w(\text{prefix}(j))$  and  $\text{target}$  is the lowest proper ancestor of node  $j$  marked with  $d$  (a dummy parent of the root node is added, marked with all the documents). See Fig. 1 for an illustration. Since the number of links with document  $\text{doc} = d$  is at most  $2|d|$ , the total number of links is at most  $\sum_{i=1}^D 2|d_i| \leq 2n$ . The following is a crucial observation by Hon et al. [12].

**Lemma 1.** *For each document  $d$  that contains a pattern  $P$ , there is a unique link with origin in the subtree of  $\text{locus}(P)$ , a proper ancestor of  $\text{locus}(P)$  as its target, and weight  $w(P, d)$ .*

We say that a link is *stabbed* by a node  $j$  if its origin is in the subtree of  $j$  ( $j$  itself included) and its target is a proper ancestor of  $j$ . Therefore, the problem of finding the  $k$ th most relevant document for  $P$  can be reduced to finding the  $k$ th highest weighted link stabbed by  $\text{locus}(P)$ .

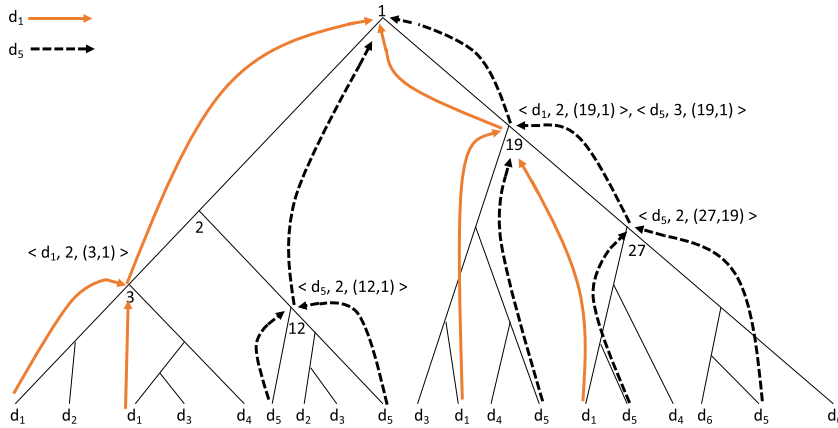


Fig. 1. An illustration of GST and links in the top- $k$  framework [12].

### 3. Super-linear space structure

In this section we start by introducing a basic data structure that uses  $O(n \log n)$  words and answers queries in  $O(\log n)$  time. Then we enhance it to a structure that uses  $O(n \log^{1+\epsilon} n)$  words, for any constant  $\epsilon > 0$ , and  $O(\log n / \log \log n)$  time. The basic structure will be used in Section 4 to achieve linear space within the same time, whereas the enhanced one will be reduced to  $O(n \log^\epsilon n)$  words. In Section 5 we show how the linear-space structure can be improved to answer queries in time  $O(\log k)$  and the enhanced structure in time  $O(\log k / \log \log n)$ , thus reaching our final results.

#### 3.1. The basic structure

We prove the following result:

**Lemma 2.** *Given the GST of a text collection of total length  $n$ , we can build an  $O(n \log n)$ -word structure that, given  $\text{locus}(P)$  and  $k$ , answers the document selection query in time  $O(\log n)$ .*

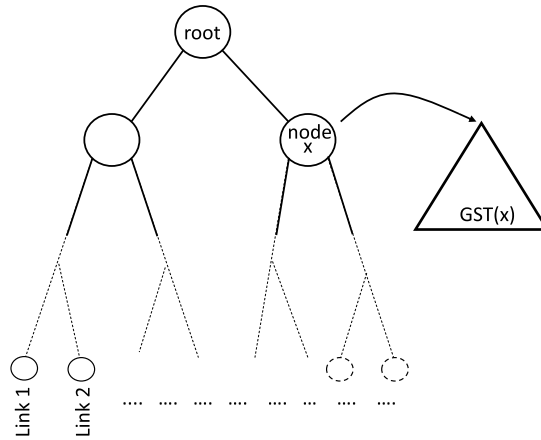
Let  $N$  represent the set of nodes in GST and  $S$  represent the set of links  $\langle \text{doc}, \text{weight}, (\text{origin}, \text{target}) \rangle$  in GST, as described in Section 2. Next we construct a balanced binary tree  $\mathcal{T}$  of  $|S|$  leaves, so that the  $i$ th highest weighted link (ties broken arbitrarily) is associated with the  $i$ th leftmost leaf of  $\mathcal{T}$ . Notice that  $n \leq |S| \leq 2n$ . We use  $S(x)$  to denote the set of links associated with the leaves in the subtree of node  $x \in \mathcal{T}$ . The intuition behind our algorithm is that if we can quickly count the number of links in any  $S(x)$  stabbed by  $\text{locus}(P)$ , then the  $k$ th smallest weighted link among all links in  $S$  stabbed by  $\text{locus}(P)$  can be easily computed via a root to leaf traversal of  $\mathcal{T}$  using  $O(\log n)$  number of such counting queries. We now present the details.

Let  $N(x)$  denote the set of nodes in GST that are (i) either the origin or the target of a link in  $S(x)$ , or (ii) the LCA of two such nodes. Clearly  $|N(x)| = \Theta(|S(x)|) = \Theta(n/2^{\text{depth}(x)})$ , where  $\text{depth}(x)$  is the number of ancestors of  $x$  (depth of root is 0). With every node  $x \in \mathcal{T}$ , we associate a tree structure  $\text{GST}(x)$ , which is the subtree of GST obtained by retaining only the nodes in  $N(x)$ , so that node  $v$  is the parent of node  $w$  in  $\text{GST}(x)$  iff  $v$  is the lowest proper ancestor of  $w$  in GST that also belongs to  $N(x)$ . See Fig. 2 for an illustration. The number of nodes and edges in  $\text{GST}(x)$  is  $\Theta(n/2^{\text{depth}(x)})$ . Notice that the same node  $w \in \text{GST}$  may appear in several  $\text{GST}(\cdot)$ 's. With each node  $w \in \text{GST}(x)$  we associate the following information:

- $\text{stab.count}_x(w)$ : The number of links in  $S(x)$  that are stabbed by node  $w$ .
- $\text{left.ptr}_x(w)$ : Let  $x_L$  be the left child of  $x$  (in  $\mathcal{T}$ ). Let  $w_L$  be the highest node in the subtree of  $w$  (in  $\text{GST}(x)$ ) that appears also in  $\text{GST}(x_L)$  ( $w_L$  can be  $w$  itself). Then  $\text{left.ptr}_x(w)$  is a pointer from  $w \in \text{GST}(x)$  to  $w_L \in \text{GST}(x_L)$ . If there exists no such node  $w_L$ , then  $\text{left.ptr}_x(w)$  is null.
- $\text{right.ptr}_x(w)$ : Analogous to  $\text{left.ptr}_x(w)$ , now considering  $x_R$ , the right child of  $x \in \mathcal{T}$ , and  $w_R$  being the highest node in the subtree of  $w \in \text{GST}(x)$  that appears also in  $\text{GST}(x_R)$ .

Note that the space needed for maintaining  $\text{GST}(x)$  and the associated information is  $O(n/2^{\text{depth}(x)})$  words. Added over all the nodes  $x \in \mathcal{T}$ , the total space occupancy of all  $\text{GST}(\cdot)$ 's is  $O(n \log n)$  words. Finally, the following result is crucial for our data structure (the case of  $w_R$  and  $x_R$  is analogous).

**Lemma 3.** *Both  $w$  and  $w_L$  stab the same subset of links of  $S(x_L)$ .*



**Fig. 2.** An illustration of links arranged in the form of a balanced binary search tree. Here the  $i$ th leftmost leaf is associated with the  $i$ th smallest weighted link. Node  $x$  is associated with  $GST(x)$ , the subtree of  $GST$  induced by the links associated with the leaves under  $x$ .

**Proof.** Otherwise, the target of a link in  $S(x_L)$  stabbing  $w_L$  but not  $w$  would be higher than  $w_L$ , below  $w$ , and belong to  $GST(x_L)$ , contradicting the definition of  $w_L$ . The same happens with the source of a link stabbing  $w$  but not  $w_L$ .  $\square$

### 3.2. Query algorithm for document selection

Assume  $locus(P)$  is given. Notice that the tree  $GST(root)$  associated with the  $root$  of  $\mathcal{T}$  is the same  $GST$  of the collection. Therefore,  $stab.count_{root}(locus(P))$  gives the number of documents containing  $P$ . If the count is less than  $k$ , there is no  $k$ th document to select. Otherwise, let  $L^*$  be the  $k$ th highest weighted link stabbed by  $locus(P)$ . Our query algorithm traverses  $\mathcal{T}$  top-down, starting from  $root$  and ending at the leaf node associated with link  $L^*$ . Then it reports the document  $d^*$  corresponding to  $L^*$ .

In our query algorithm, we use  $x$  to denote a node in  $\mathcal{T}$ ,  $w$  to denote a node in  $GST(x)$  and  $K$  to denote an integer  $\leq k$ . First we initialize  $x$  to the root of  $\mathcal{T}$ ,  $w$  to  $locus(P)$  and  $K$  to  $k$ . This establishes the invariant that we have to return the  $K$ th highest weighted link in  $S(x)$  stabbed by  $w$ . Let  $x_L$  and  $x_R$  be the left and right children of  $x$ . Then we obtain the nodes  $w_L \in GST(x_L)$  and  $w_R \in GST(x_R)$  pointed by  $left.ptr_x(w)$  and  $right.ptr_x(w)$ , respectively. The following values are then computed in constant time.

- $c = stab.count_x(w)$ , the number of links in  $S(x)$  stabbed by  $w$ .
- $c_L = stab.count_{x_L}(w_L)$ , the number of links in  $S(x_L)$  stabbed by  $w$  (or  $w_L$ ).
- $c_R = stab.count_{x_R}(w_R)$ , the number of links in  $S(x_R)$  stabbed by  $w$  (or  $w_R$ ).

Notice that  $c = c_L + c_R$ . If  $c_L \geq K$  then, by Lemma 3, the  $K$ th link below  $S(x)$  (or  $S(x_L)$ ) stabbed by  $w \in GST(x)$  is the same as the  $K$ th link below  $S(x_L)$  stabbed by  $w_L \in GST(x_L)$ . Therefore, we maintain the invariant if we continue the traversal in the subtree of  $x \leftarrow x_L$  with  $GST(x_L)$  node  $w \leftarrow w_L$ . On the other hand, if  $c_L < K$ , then by Lemma 3 the  $K$ th link stabbed by  $w$  below  $S(x)$  is the same as the  $(K - c_L)$ th link below  $S(x_R)$  stabbed by  $w_R \in GST(x_R)$ . In this case, we maintain the invariant if we continue the traversal in the subtree of  $x \leftarrow x_R$  with  $GST(x_R)$  node  $w \leftarrow w_R$  and with  $K \leftarrow K - c_L$ . We terminate the algorithm when  $x$  is a leaf, thus  $K = 1$  and  $x$  represents  $L^*$ . As the height of  $\mathcal{T}$  is  $O(\log n)$  and the time spent at each node is constant, the total query time is  $O(\log n)$  and Lemma 2 is proved.

### 3.3. An enhanced structure

We now prove the following result, which holds in the RAM model of computation, with a computer word of  $w = \Omega(\log n)$  bits.

**Lemma 4.** Given the  $GST$  of a text collection of total length  $n$  and any constant  $0 < \epsilon \leq 1$ , we can build an  $O(n \log^{1+\epsilon} n)$ -word structure that, given  $locus(P)$  and  $k$ , answers the document selection query in time  $O(\log n / \log \log n)$ .

In order to speed up the structure of Lemma 2, we will choose a step  $s = \epsilon \log \log n$  and build the  $GST(x)$  structures only for nodes  $x \in \mathcal{T}$  whose depth is a multiple of  $s$ . Each node  $w \in GST(x)$  for the selected nodes  $x$  will store sufficient information for the query algorithm to jump directly to the corresponding node  $x'$  at depth  $depth(x') = depth(x) + s$ , instead of just to  $x_L$  or  $x_R$ .

Given  $x, x' \in \mathcal{T}$  as above ( $x'$  in the subtree of  $x$ ) and  $w \in \text{GST}(x)$ , we define  $w_{x'}$  as the highest node in the subtree of  $w$  that appears also in  $\text{GST}(x')$ . Let us call  $x_1, x_2, \dots, x_{2^s}$  the nodes at depth  $\text{depth}(x) + s$  that descend from  $x$  (or the leaves below  $x$ , if they have depth less than  $\text{depth}(x) + s$ ), ordered left to right in  $\mathcal{T}$  (i.e., from highest to lowest weights in  $S(x_i)$ ).

Associated to each node  $w \in \text{GST}(x)$ , we store  $2^s$  pointers  $\text{ptr}_x(w)[i] = w_{x_i}$ . We also store the  $2^s$  cumulative values  $\text{acc}_x(w)[i] = \sum_{j=1}^i \text{stab.count}_{x_j}(w_{x_j})$ ; note that  $\text{acc}_x(w)[2^s] = \text{stab.count}_x(w)$ . We will store those  $\text{acc}_x(w)$  values in a fusion tree [7], which takes  $O(2^s) = O(\log^\epsilon n)$  words of space and solves predecessor queries in  $\text{acc}_x(w)$  in constant time. The space is the same used by array  $\text{ptr}_x(w)$ , which added over all the  $\text{GST}(\cdot)$ 's is  $O(n \log^{1+\epsilon} n)$  words (even if only one level out of  $s$  in  $\mathcal{T}$  stores  $\text{GST}(\cdot)$  structures).

Queries now proceed as in Section 3.2, but now we use the fusion tree to determine, given  $w \in \text{GST}(x)$ , which is the node  $x_i \in \mathcal{T}$  that contains the  $K$ th link below  $S(x)$  stabbed by  $w$ . Therefore we can move directly from  $x$  to  $x_i$  and from  $w \in \text{GST}(x)$  to  $w_i \in \text{GST}(x_i)$ , where  $w_i = \text{ptr}_x(w)[i]$ . We also update  $K \leftarrow K - \text{acc}_x(w)[i - 1]$  (assume  $\text{acc}_x(w)[0] = 0$ ). Thus we complete the query in  $O((\log n)/s) = O(\log n/(\epsilon \log \log n))$  constant-time steps and Lemma 4 is proved.

#### 4. Linear space structure

In this section we build on the basic structure of Lemma 2 in order to achieve linear space and logarithmic query time. At the end, we reduce the space of the enhanced structure to  $O(n \log^\epsilon n)$ . The results hold under the RAM model.

**Lemma 5.** *Given the GST of a text collection of total length  $n$ , we can build an  $O(n)$ -word structure that, given  $\text{locus}(P)$  and  $k$ , answers the document selection query in time  $O(\log n)$ .*

To achieve linear space, we replace some of our data structures by succinct ones. We will measure the space in bits, aiming at using  $O(n \log n)$  bits overall. The binary tree  $\mathcal{T}$  can be maintained in  $O(n \log n)$  bits, where each internal node  $x$  stores an  $O(\log n)$ -bit pointer to the corresponding tree  $\text{GST}(x)$  and each leaf stores the document identifier corresponding to the associated link. The global GST can also be maintained in  $O(n \log n)$  bits. Therefore, the space-consuming components are the  $\text{GST}(\cdot)$ 's and their associated information.

Using well-known succinct data structures [26], the  $\text{GST}(x)$  tree topologies can be represented in  $O(1)$  bits per node (i.e.,  $O(n \log n)$  bits overall) with constant-time support of all the basic navigational operations required in our algorithm. We refer to any node  $w \in \text{GST}(x)$  by its pre-order rank, that is, node  $j$  means the node with pre-order rank  $j$ . The pre-order rank of the root node of any  $\text{GST}(x)$  is 1. Next we show how to encode the remaining information associated with each node in  $\text{GST}(x)$  using  $O(1)$  bits per node.

##### 4.1. Encoding $\text{stab.count}_x(j)$

We note that  $\text{stab.count}_x(j)$  is exactly equal to the number of links of  $S(x)$  associated with  $\text{GST}(x)$  that originate in the subtree of  $j$  minus the number of links in  $S(x)$  that target any node in the subtree of  $j$  ( $j$  belongs to its subtree). We encode this information in two bit vectors:  $B_x = 10^{\alpha_1} 10^{\alpha_2} 10^{\alpha_3} \dots$  and  $B'_x = 10^{\beta_1} 10^{\beta_2} 10^{\beta_3} \dots$ , where  $\alpha_j$  (resp.,  $\beta_j$ ) is the number of links of  $S(x)$  originating from (resp., targeting at) node  $j$  in  $\text{GST}(x)$ . We augment  $B_x$  and  $B'_x$  with structures supporting constant-time rank/select queries [18]. Notice that  $\sum \alpha_j = \sum \beta_j = O(|S(x)|) = O(|\text{GST}(x)|)$ . Therefore, both  $B_x$  and  $B'_x$  can be represented in  $O(1)$  bits per node.

Now we can compute  $\text{stab.count}_x(j)$  for any  $j$  in  $O(1)$  time as follows: find the rightmost leaf node  $j'$  in the subtree of  $j$  in  $O(1)$  time using the succinct tree representation of  $\text{GST}(x)$  [26]. Then the number  $n_o$  of links originating from the subtree of  $j$  is equal to the number of 0-bits between the  $j$ th and  $(j' + 1)$ th 1-bit in  $B_x$  (because  $j$  and  $j'$  are preorder numbers). Similarly, the number  $n_t$  of links targeted at any node in the subtree of  $j$  is equal to the number of 0-bits between the  $j$ th and  $(j' + 1)$ th 1-bits in  $B'_x$ . Using rank/select operations on  $B_x$  and  $B'_x$ ,  $n_o$  and  $n_t$  are computed in  $O(1)$  time and  $\text{stab.count}_x(j)$  is given by  $n_o - n_t$ .

##### 4.2. Encoding $\text{left.ptr}_x(j)$ and $\text{right.ptr}_x(j)$

We show how to encode  $\text{left.ptr}_x(\cdot)$  for all nodes in  $\text{GST}(x)$ ;  $\text{right.ptr}_x(j)$  is symmetric. The idea is to maintain a bit vector  $LP$  such that  $LP[j] = 1$  iff there exists a node  $j_L \in \text{GST}(x_L)$  such that both  $j \in \text{GST}(x)$  and  $j_L \in \text{GST}(x_L)$  represent the same node in GST. We add constant-time rank/select data structures [18] on  $LP$ . Since the length of  $LP$  is equal to the number of nodes in  $\text{GST}(x)$ , its space occupancy is  $O(1)$  bits per node.

Now, for any given node  $j \in \text{GST}(x)$ , the node  $j_L \in \text{GST}(x_L)$  to which  $\text{left.ptr}_x(j)$  points is the (unique) highest descendant of  $j$  that is marked in  $LP$ , thus it can be identified by (1) finding the position  $j^*$  of the leftmost 1-bit in  $LP[j \dots]$ ; (2) checking if node  $j^*$  is in the subtree of node  $j$  in  $\text{GST}(x)$ ; (3) if so, then  $j_L \in \text{GST}(x_L)$  is equal to the number of 1's in  $LP[1 \dots j^*]$ , otherwise,  $j_L$  is null. All these operations require constant time, either using the succinct tree operations or the rank/select data structures. This works because all the nodes in  $\text{GST}(x_L)$  appear in  $\text{GST}(x)$ , in the same order (pre-order).

In summary, the space requirement of our encoding scheme is  $O(1)$  bits per node in any  $\text{GST}(x)$ , thus adding to  $O(n \log n)$  bits. The query algorithm, as well as its time complexity, remain the same. This completes the proof of Lemma 5.

### 4.3. Reducing space of the enhanced structure

The space of the enhanced structure of Section 3.3 can be similarly reduced to  $O(n \log^\epsilon n)$  words, obtaining the following result.

**Lemma 6.** *Given the GST of a text collection of total length  $n$  and a constant  $\epsilon > 0$ , we can build an  $O(n \log^\epsilon n)$ -word structure that, given  $\text{locus}(P)$  and  $k$ , answers the document selection query in time  $O(\log n / \log \log n)$ .*

For this sake, recalling the definition of  $x_1, \dots, x_{2^s}$  of Section 3.3, we will maintain bit vectors  $LP_i$  for  $i = 1$  to  $2^s$ , so that  $LP_i[j] = 1$  iff there exists a node  $j_i \in \text{GST}(x_i)$  such that both  $j \in \text{GST}(x)$  and  $j_i \in \text{GST}(x_i)$  represent the same node in GST. Then each array entry  $\text{ptr}_x(j)[i]$  is computed using  $LP_i$  as in Section 4.2. The total space used by all the  $LP_i$  bit vectors is  $O(2^s) = O(\log^\epsilon n)$  bits per node, adding up to  $O(n \log^{1+\epsilon} n)$  bits in total.

To compute  $\text{acc}_x(j)[i]$ , we store bitmaps  $B_{x,1}, \dots, B_{x,2^s}$  and  $B'_{x,1}, \dots, B'_{x,2^s}$ , analogous to  $B$  and  $B'$  of Section 4.1. In this case,  $B_{x,i} = 10^{\alpha_1} 10^{\alpha_2} 10^{\alpha_3} \dots$ , so that  $\alpha_j^i = \sum_{r=1}^i s(r)$ , where  $s(r)$  is the number of links of  $S(x_r)$  originating from node  $\text{ptr}_x(j)[i] \in \text{GST}(x_r)$ , and  $B'_{x,i} = 10^{\beta_1} 10^{\beta_2} 10^{\beta_3} \dots$ , so that  $\beta_j^i = \sum_{r=1}^i t(r)$ , where  $t(r)$  is the number of links of  $S(x_r)$  targeting at node  $\text{ptr}_x(j)[i] \in \text{GST}(x_r)$ . Then, it holds  $\text{acc}_x(j)[i] = \alpha_j^i - \beta_j^i$ , which is computed in constant time using rank/select operations. Since it holds  $\alpha_j^i \leq \alpha_j$  and  $\beta_j^i \leq \beta_j$  for all  $i$  values, the total space of these  $2^s = \log^\epsilon n$  bitmaps adds up to  $O(n \log^{1+\epsilon} n)$  bits.

To carry out predecessor searches on the virtual vector  $\text{acc}_x(j)$ , we use succinct SB-trees [8, Lemma 3.3]. Given constant-time access to any  $\text{acc}_x(j)[i]$ , this structure provides predecessor searches in  $O(1 + \log(2^s) / \log \log n) = O(1)$  time and use  $O(2^s \log \log n) = O(\log^\epsilon n)$  bits per node (by adjusting  $\epsilon$ ). Thus the total space is  $O(n \log^{1+\epsilon} n)$  bits as well. This concludes the proof of Lemma 6.

## 5. Achieving $O(\log k)$ query time and better

In this section we first build on the linear-space data structure of Lemma 5 in order to improve its query time to  $O(\log k)$ . At the end, we show that the result extends to our superlinear-space data structure of Lemma 6, improving its query time to  $O(\log k / \log \log n)$ . Thus we start by proving the following theorem.

**Theorem 1.** *A collection  $\mathcal{D}$  of documents can be preprocessed into a linear-space data structure that can answer any document selection query  $(P, k)$  in time  $O(\log k)$ , given the locus of pattern  $P$  in the generalized suffix tree of  $\mathcal{D}$ .*

Notice that the query time  $O(\log n)$  in Lemma 5 can be written as  $O(\log k)$  for  $k > \sqrt{n}$ . Therefore, we turn our attention to the case where  $k \leq \sqrt{n}$ . First, we derive a space-efficient structure  $DS(\delta)$ , which can answer document selection queries faster, but only for values of  $k$  below a predefined parameter  $\delta \leq \sqrt{n}$ . More precisely, structure  $DS(\delta)$  will satisfy the following properties:

**Lemma 7.** *The structure  $DS(\delta)$  uses  $O(n(\log \delta + \log \log n))$  bits of space and can answer document selection queries in time  $O(\log \delta + \log \log n)$ , for  $k \leq \delta \leq \sqrt{n}$ .*

To obtain the result in Theorem 1, we maintain structures  $DS(\delta_i)$  with  $\delta_i = \lceil n^{1/2^i} \rceil$  for  $i = 1, 2, 3, \dots, r$ , where  $\delta_{r+1} \leq \sqrt{\log n} < \delta_r$  (therefore  $r < \log \log n$ ). The total space needed is  $O(n \sum_{i=1}^r (\log \delta_i + \log \log n)) = O(n \log n)$  bits ( $O(n)$  words). When  $k$  comes as a query, if  $k > \delta_{r+1}$ , we first find  $h$ , where  $\delta_{h+1} < k \leq \delta_h$  and obtain the answer using  $DS(\delta_h)$ . The resulting time is  $O(\log \delta_h + \log \log n) = O(\log k)$ . The case where  $k < \delta_{r+1}$  is handled separately using other structures in  $O(1)$  time (Section 5.2). We now describe the details of  $DS(\delta)$ .

### 5.1. Structure $DS(\delta)$

The first step is to identify certain nodes in GST as *marked* nodes and *prime* nodes, based on a parameter  $g = \lceil \delta \log n \rceil$  called the *grouping factor*. Every  $g$ th leftmost leaf is marked, and the LCA of every two consecutive marked leaves is also marked. Therefore, the number of marked nodes is  $\Theta(n/g)$ . Nodes with their parent marked are prime. A prime node with at least one marked node in its subtree is a type-1 prime node, otherwise it is a type-2 prime node. Notice that the highest marked node in the subtree of any node is unique, if it exists. Therefore, except the root node, every marked node  $j^*$  can be associated with a unique type-1 prime node  $j'$ , which is the first prime node on the path from  $j^*$  to the root. Notice that a node can be both prime and marked.

Let  $j'$  be a prime node and  $j^*$  be the highest marked node in its subtree ( $j^*$  exists only if  $j'$  is of type-1, and it can be that  $j' = j^*$ ). We use  $G(j' \setminus j^*)$  to represent the subtree of GST rooted at  $j'$  after removing the subtree of  $j^*$  ( $j^*$  is not removed). With a slight abuse of notation, we use  $G(j' \setminus j^*)$  to represent the set of nodes within  $G(j' \setminus j^*)$  as well. A crucial result [32] is that, for any prime node  $j'$ , the number of nodes in  $G(j' \setminus j^*)$  is  $O(g)$ .

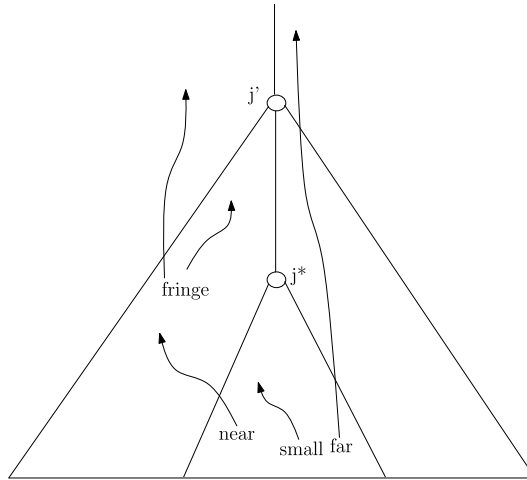


Fig. 3. Categorization of Links [32].

We define  $prime.parent(j)$  of any node  $j$  in GST as the first prime node  $j'$  on the path from  $j$  to the root. Note that  $j \in G(j' \setminus j^*)$ , otherwise  $j$  would be a (strict) descendant of  $j^*$  and its corresponding  $j'$  would be below  $j^*$ .

It is not hard to determine  $j' = prime.parent(j)$  in constant time and  $O(n)$  bits, by sampling the prime nodes in a succinct tree representation and looking for the lowest sampled ancestor of  $j$  [31, Lemma 4.4].

The structure  $DS(\delta)$  is a collection of substructures  $STR(j')$  associated with every prime node  $j'$  in GST. If the input node  $locus(P) \in G(j' \setminus j^*)$  and  $k \leq \delta$ , we obtain the answer using  $STR(j')$  in  $O(\log g) = O(\log \delta + \log \log n)$  time. Based on the type of  $j'$ , we have two cases; we describe the simpler one first.

### 5.1.1. $STR(j')$ associated with a type-2 prime node $j'$

The structure can be constructed as follows: take  $G(j')$ , the subtree rooted at node  $j'$ , and replace the pre-order rank of each node  $j$  by  $(j - j' + 1)$ . Also associate a dummy parent node to the root. Then, among the links defined over GST (Section 2), choose those that originate from the subtree of  $j'$  and: (1) Assign a new value to its origin and target, which is its original value minus  $j'$  plus 1. The target of some links can be negative; replace those by 0. (2) Replace the weight by a rank-space reduced value in  $[1, O(|G(j')|)]$ . Notice that the number of links chosen is  $O(|G(j')|)$ . (3) Let  $d$  be its document identifier. Instead of writing  $d$  explicitly in  $\lceil \log D \rceil$  bits, use a pointer to one leaf node in  $G(j')$ , using  $\lceil \log |G(j')| \rceil$  bits, where the suffix corresponding to that leaf belongs to document  $d$ .

In summary, we have a tree of  $(|G(j')| + 1)$  nodes and  $O(|G(j')|)$  links associated with it. The information (origin, target, document, weight) associated with each link is encoded in  $O(\log |G(j')|)$  bits. Then  $STR(j')$  is the structure described in Lemma 5 over these nodes and links. The space required is  $O(|G(j')| \log |G(j')|) = O(|G(j')| \log g)$  bits. We maintain structures  $STR(j')$  for all type-2 prime nodes  $j'$  in total  $O(n \log g)$  bits, since a node can be in the subtree of at most one type-2 prime node.

### 5.1.2. $STR(j')$ associated with a type-1 prime node $j'$

We first identify the candidate set  $\mathcal{C}(j')$  of  $O(g)$  links, such that for any  $k \leq \delta$ , the  $k$ th link stabbed by any node  $j \in G(j' \setminus j^*)$  belongs to  $\mathcal{C}(j')$ . Clearly we can ignore the links that do not originate from the subtree of  $j'$ . The links that do can be categorized into the following types [32]: near-links are stabbed by  $j^*$ , but not by  $j'$ ; far-links are stabbed by both  $j^*$  and  $j'$ ; small-links are targeted at a node in the subtree of  $j^*$ ; and fringe-links are the others (see Fig. 3).

We include all near-links and fringe-links into  $\mathcal{C}(j')$ , which are  $O(g)$  in number [32, Lemma 8]. All small-links can be ignored as none of them is stabbed by any node in  $G(j' \setminus j^*)$ . Notice that if any node in  $G(j' \setminus j^*)$  stabs a far-link, it indeed stabs all far-links. Therefore, it is sufficient to insert the top- $\delta$  far-links into  $\mathcal{C}(j')$ . Thus, we have  $O(g)$  links in  $\mathcal{C}(j')$  overall.

Now we perform a rank-space reduction of pre-order rank of nodes in  $G(j' \setminus j^*)$  as well as of the information associated with the links in  $\mathcal{C}(j')$ , as follows:

- The target of those links targeting at any proper ancestor of  $j'$  is changed to a dummy parent node of  $j'$ . Similarly, the origin of all those links originating in the subtree of  $j^*$  is changed to node  $j^*$ .
- The pre-order rank of all those nodes in  $G(j' \setminus j^*)$ , and the corresponding origin and target values of links in  $\mathcal{C}(j')$ , are changed to a rank-space reduced value in  $[0, |G(j' \setminus j^*)|]$ . Notice that the new pre-order rank of  $j'$  is 1 and that of its dummy parent node is 0. We remark that this mapping (and remapping) can be stored separately in  $O(|G(j' \setminus j^*)| \log |G(j' \setminus j^*)|)$  bits.
- The weights of the links are also replaced by rank-space reduced values.

- Let  $L$  be a near- or fringe-link in  $\mathcal{C}(j')$  with  $d$  its corresponding document. Then there must be at least one leaf  $\ell$  in  $G(j' \setminus j^*)$  where the suffix corresponding to  $\ell$  belongs to  $d$ . Therefore, instead of representing  $d$ , we maintain a pointer to  $\ell$ , which takes only  $O(\log g)$  bits. This trick will not work for far-links, as the existence of such a leaf node is not guaranteed. Therefore, we spend  $\log D$  bits for each far-link, which is still affordable because there are only  $O(\delta) = O(g/\log n)$  far-links.

In summary, we have a tree of  $(|G(j' \setminus j^*)| + 1) = O(g)$  nodes with  $O(g)$  links associated with it. Then  $STR(j')$  is the structure described in Lemma 5 over these nodes and links. The space required is  $O(g \log g)$  bits. As the number of type-1 prime nodes is  $O(n/g)$ , the total space to maintain  $STR(j')$  for all type-2 primes nodes  $j'$  is  $O(n \log g)$  bits.

### 5.1.3. Query answering

Given node  $j = \text{locus}(P)$ , we find  $j' = \text{prime.parent}(j)$ . Then we map node  $j$  to the corresponding node in  $STR(j')$  and obtain the answer by querying  $STR(j')$ , in  $O(\log g) = O(\log \delta + \log \log n)$  time. The answer may come in the form of a node in  $STR(j')$ , which is mapped back to GST in order to obtain the associated document. This completes the proof of Lemma 7.

### 5.2. Structure for $k \leq \delta_{r+1}$

First, identify the marked and prime nodes in GST with  $g = \delta_{r+1} \log n$ . At every prime node  $j'$ , we explicitly maintain the candidate set  $\mathcal{C}(j')$ . This takes  $O(n)$ -word space. Then for any  $k \leq \delta_{r+1}$ , the  $k$ th link stabbed by node  $j$  can be encoded as a pointer to the corresponding entry in  $\mathcal{C}(\text{prime.parent}(j'))$  using  $\lceil \log |\mathcal{C}(\text{prime.parent}(j'))| \rceil = O(\log g) = O(\log \log n)$  bits. Therefore, the answers for all  $k \in [1, \delta_{r+1}]$  for all nodes in GST can be maintained in additional  $O(n \cdot \delta_{r+1} \log \log n) = o(n \log n)$  bits of space. Now the  $k$ th link (and its document) stabbed by any query node  $\text{locus}(P)$  can be obtained from  $\mathcal{C}(\text{prime.parent}(\text{locus}(P)))$  in  $O(1)$  time.

### 5.3. Speeding Up the enhanced structure

The same construction used above can be used to speed up our superlinear-space structure of Lemma 6, simply by using it instead of the linear-space one of Lemma 5 to implement the structures  $STR(j')$ . The space  $O(n \log^\epsilon n)$  words, or  $O(n \log^{1+\epsilon} n)$  bits, will become  $O(g \log g \log^\epsilon n)$  inside the structures  $STR(j')$ , because we will maintain the sampling step  $s = \epsilon \log \log n$  depending on  $n$ , not on  $g$ , and use the succinct SB-trees with parameter  $n$ , not  $g$ . As a result, the total space per value of  $\delta$  will be  $O(n \log g \log^\epsilon n)$  bits, and added over all the values of  $\delta$  we will have  $O(n \log^\epsilon n \sum_{i=1}^r (\log \delta_i + \log \log n)) = O(n \log^{1+\epsilon} n)$  bits, or  $O(n \log^\epsilon n)$  words. The time, on the other hand, will be  $O(1 + \log \delta / (\epsilon \log \log n))$  on  $DS(\delta)$ , which becomes  $O(1 + \log k / (\epsilon \log \log n))$  in terms of  $k$ . We have proved our final result for the superlinear structure.

**Theorem 2.** *A collection  $\mathcal{D}$  of documents of total length  $n$  can be preprocessed into a data structure using  $O(n \log^\epsilon n)$  words of space, for any constant  $\epsilon > 0$ , which can answer document selection queries  $(P, k)$  in time  $O(1 + \log k / \log \log n)$ , given the locus of pattern  $P$  in the generalized suffix tree of  $\mathcal{D}$ .*

## 6. Lower bounds

In this section we prove two lower bounds. The first shows that it is highly unlikely that an efficient succinct-space solution for range selection exists. The second shows, under a somewhat weak model, that  $O(1 + \log k / \log \log n)$  is the best possible time under  $O(n \text{polylog } n)$  space.

### 6.1. Hardness of an efficient succinct solution

One could expect to obtain an index using  $O(n \log \sigma)$  bits of space, proportional to the  $n \log \sigma$  bits needed to store  $\mathcal{D}$ , as achieved for the top- $k$  document retrieval problem. We show, however, that this is very unlikely unless a significant breakthrough in the current state of the art of computational geometry is obtained, specifically on the three-dimensional orthogonal range reporting problem. The best known space time trade-offs are  $O(n \log^{1+\epsilon} n)$  space with  $O(\log \log n + \text{output})$  time, and  $O(n \log n)$  space with  $O((1 + \text{output}) \log^\epsilon n)$  time [4]. Within linear space, the best known time complexities are  $O(n^{2/3} + \text{output})$  [16] and  $O(n^{1/3} + \text{output}) \log n$  [30].

**Theorem 3.** *If there exists a data structure using  $O(n \log \sigma + D \text{polylog } n)$  bits with query time  $O(|P| \text{polylog } n)$  for document selection ( $\sigma$  being the alphabet size), then there exists a linear-space data structure that can answer three-dimensional range reporting queries in poly-logarithmic time per query plus per reported point.*

**Proof.** We reduce from the position restricted substring searching (PRSS) problem, which is defined as follows: Index a given a text  $T[1, n]$  over an alphabet set  $[1, \sigma]$ , such that whenever a pattern  $P$  (of length  $p$ ) and a range  $[x, y]$  comes as a



query, all those  $occ_{x,y}$  occurrences of  $P$  in  $T[x \dots y]$  can be reported efficiently. Many indexes offering different space and query time trade-offs exist [2,17,29].

Hon et al. [11] proved that answering PRSS queries in polylog time and succinct space is at least as hard as performing 3-dimensional orthogonal range reporting in polylog time and linear space. They also showed that if the query pattern is longer than  $\alpha = \lceil \log^{2+\epsilon} n \rceil$  for some predefined constant  $\epsilon > 0$ , an efficient succinct space index can be designed. Therefore, the harder case arises when  $p < \alpha$ . We now show how to answer PRSS queries with  $p < \alpha$  via document selection queries on the following set:  $\mathcal{D} = \{d_1, d_2, d_3, \dots, d_{\lceil n/\alpha \rceil}\}$ , where  $d_i = T[1 + (i - 1)\alpha \dots (i + 1)\alpha]$  and  $|d_i| = 2\alpha$ , except possibly for  $d_{\lceil n/\alpha \rceil - 1}$  and  $d_{\lceil n/\alpha \rceil}$ . The score function  $w(P, d_i)$  is  $i$  if  $P$  appears at least once in  $d_i$  and 0 otherwise. Notice that an occurrence of any pattern of length at most  $\alpha$  overlaps with at least one and at most two documents in  $\mathcal{D}$ . Therefore, the previously defined PRSS query on  $T$  can be answered via multiple document selection queries on  $\mathcal{D}$  as follows: first report all those documents  $d_i$  with  $w(P, d_i) \in [\lceil x/\alpha \rceil, \lfloor y/\alpha \rfloor + 2]$ . Then, within all those reported documents, look for other occurrences of  $P$  via an exhaustive scanning. If the time for document selection queries is polylog in the total length of all documents in  $\mathcal{D}$  (which is at most  $2n$ ), then the time for PRSS query is also bounded by  $O((p + occ_{x,y}) \text{polylog } n)$ . Therefore, answering document selection queries in polylog time and succinct space is at least as hard as answering PRSS queries in polylog time and succinct space.  $\square$

### 6.2. A weak lower bound for range selection time

In this section we show that  $O(1 + \log k / \log \log n)$  is the best possible time unless we use more than  $O(n' \text{polylog } n)$  space, where  $n'$  is the number of internal nodes in GST. In most text families,  $n'$  will be  $\Theta(n)$ , but there are cases, for example on very repetitive text collections, where it may very well be that  $n' = o(n)$ .

**Theorem 4.** Any data structure using  $O(n' \text{polylog } n)$  words of space to solve the document selection problem must have query time  $\Omega(1 + \log k / \log \log n)$ , where  $n$  is the collection size and  $n'$  is the number of internal nodes in its GST.

**Proof.** We reduce from the prefix selection problem [15], which has this lower-bound complexity within  $O(n \text{polylog } n)$  space. In this problem we have an array  $A[1, n]$  and queries ask for the  $k$ th largest element in some prefix  $A[1, i]$ , for any  $k$  and  $i$ . The query time of  $O(1 + \log k / \log \log n)$  is achievable using linear space [5].

Now consider the document set formed by  $D = n$  documents where document  $d_i = a^{n-i}b$ . The pattern  $P_i = a^{n-i}b$  appears in documents  $d_1$  to  $d_i$ . We will use  $O(n)$  space to precompute the locus nodes of all patterns  $P_i$  in GST. Finally, we assign weight  $A[i]$  to the document  $d_i$ .

Then,  $A[j]$  is the  $k$ th largest element in  $A[1, i]$  iff  $d_j$  is the  $k$ th heaviest document among  $d_1, \dots, d_i$ , that is, among the documents where pattern  $P_i$  appears. The document selection query gives the answer to the range selection query.

Note that the total length of the documents is  $\Theta(n^2)$ , thus the  $n$  values are different in both problems. However, the GST of this particular collection has  $n' = \Theta(n)$  internal nodes. Note  $\log |\text{GST}| = \Theta(\log n)$ .

Thus, if we could solve document selection queries in time  $o(\log k / \log \log n)$  within  $O(n' \text{polylog } n)$  words of space in any text collection, we would have a tool to solve any arbitrary prefix selection query in  $A[1, n]$  in  $o(\log k / \log \log n)$  time and within  $O(n \text{polylog } n)$  words of space, a contradiction.  $\square$

## 7. Further space-time tradeoffs

In this section, we present some additional results that complement the lower bounds of the previous sections. As for the space lower bound, we show that  $O(n)$  bits of space are sufficient if we accept a small error in the ranking. As for the time, we show that constant-time is possible if we spend space  $O(n^{1+\epsilon})$  for any constant  $\epsilon > 0$ .

### 7.1. Efficient approximate term-frequency based retrieval

In this section, we show how to obtain efficient data structures, both in terms of space as well as the query time, for the case where the ranking function  $w(P, d) = \text{TF}^*(P, d)$ , where  $\text{TF}^*(P, d)$  is an approximation of term-frequency  $\text{TF}(P, d)$  (the number of occurrences of  $P$  in document  $d$ ). Specifically,

$$0 \leq \text{TF}^*(P, d) - \text{TF}(P, d) \leq 2\epsilon \log n, \text{ for any constant } \epsilon > 0$$

Our main result is summarized below.

**Theorem 5.** A collection  $\mathcal{D}$  of documents of total length  $n$  can be preprocessed into a data structure using  $O(n)$  bits of space, which can answer the  $k$ th most relevant document for any pattern  $P$  w.r.t. the ranking function  $\text{TF}^*(P, \cdot)$  in time  $(\log k)$ , given the suffix range of pattern  $P$  in a generalized suffix tree of  $\mathcal{D}$ .

**Proof.** Our idea is to first construct a sampled generalized suffix tree of  $\mathcal{D}$  (denoted by  $\text{GST}^\epsilon$ ) as follows:  $\text{GST}^\epsilon$  is a combined compact trie of every  $(\epsilon \log n)$ th lexicographically smallest suffixes of all the documents in  $\mathcal{D}$ . Therefore  $\text{GST}^\epsilon$  will be

a trie of  $(n/\log n)$  nodes. We then apply our linear space framework over  $\text{GST}^s$  and obtain an  $O(n/\log n)$  words (i.e.,  $O(n)$  bits) index. We also maintain a bitmap  $B[1..n]$  such that  $B[i] = 1$  iff the  $i$ th lexicographically smallest suffix in  $\text{GST}$  belongs to  $\text{GST}^s$ . Structures supporting rank/select queries over  $B$  in constant time is also maintained, in additional  $o(n)$  bits.

To find the  $k$ th link stabbed by  $\text{locus}(P)$ , we first identify the suffix range  $[sp, ep]$  of  $P$ . Then identify  $sp'$  and  $ep'$ , where  $sp'$  is 1 plus the number of 1's in  $B[1..sp - 1]$  and  $ep'$  is the number of 1's in  $B[1..ep]$ . This step requires only constant time. Now, the  $k$ th index stabbed by the LCA of  $sp'$ 'th and  $ep'$ 'th leaves in  $\text{GST}^s$  is our answer and can be computed in  $O(\log k)$  time using our  $O(n)$  bit index. This works because the number of suffixes of any document  $d$  in  $\text{GST}^s$  prefixed by  $P$  is at most  $\text{TF}(P, d)/(\epsilon \log n)$  and at least  $\text{TF}(P, d)/(\epsilon \log n) - 2$ .  $\square$

The following result can be obtained by applying the framework by Shah et al. [32] over  $\text{GST}^s$ .

**Theorem 6.** *A collection  $\mathcal{D}$  of documents of total length  $n$  can be preprocessed into a data structure using  $O(n)$  bits of space, which can retrieve the top- $k$  most relevant documents for any pattern  $P$  w.r.t. the ranking function  $\text{TF}^*(P, \cdot)$  in time  $O(k)$ , given the suffix range of pattern  $P$  in a generalized suffix tree of  $\mathcal{D}$ .*

## 7.2. A constant-query-time index

We present the following result in this section.

**Theorem 7.** *A collection  $\mathcal{D}$  of documents of total length  $n$  can be preprocessed into a data structure using  $O(n^{1+\epsilon})$  words of space, for any constant  $\epsilon > 0$ , which can answer document selection queries  $(P, k)$  in constant time, given the locus of pattern  $P$  in the generalized suffix tree of  $\mathcal{D}$ .*

**Proof.** The result can be obtained by several simple modifications of our super-linear space structure described in Section 3. First, increase the fan out of the basic tree structure to  $\Theta(n^\epsilon)$  instead of limiting by 2. Thus the tree height becomes a constant and the overall space for the associated information will be linear. Before introducing an additional data structure, we briefly describe how the query algorithm changes. It is again a root to leaf traversal of  $\mathcal{T}$ . Let  $x$  be a node in  $\mathcal{T}$  and  $w$  be a node in  $\text{GST}(x)$ . Then to find the  $K$ th link stabbed by  $w$ , we do the following. For  $i = 1, 2, 3, \dots, \text{degree}(x)$ , find

- $x_i$ , the  $i$ th leftmost child of  $x$
- $w_i$ , the node in  $\text{GST}(x_i)$  corresponding to  $w$  in  $\text{GST}(x)$
- $\text{stab.count}_{x_i}(w_i)$

Next we compute the smallest  $c$  such that  $\sum_{i=1}^c \text{stab.count}_{x_i}(w_i) \geq K$ . Then, instead of asking for the  $K$ th link stabbed by  $w$ , it is equivalent to ask for the  $(\sum_{i=1}^c \text{stab.count}_{x_i}(w_i) - K)$ th link stabbed by  $w_c$ . However, the above steps require  $O(\text{degree}(x)) = O(n^\epsilon)$  time, and so does the total query time.

We now show how to achieve constant query time by associating the following information with every node  $w$  (as described above). Store an array  $A_w[1..n^\epsilon]$ , such that  $A_w[i] = \text{stab.count}_x(w)$ . Also maintain another array  $A'_w$ , where  $A'_w[i] = \sum_{j=1}^i A_w[j]$ . We then maintain a predecessor search structure over  $A'_w$  with constant query time. Notice that there exists such a structure, requiring quadratic space [6]. Therefore by associating an  $O(n^{2\epsilon})$  space structure with  $w$ , the step where we identify  $w_c$  can be performed in constant time via a single predecessor search query on  $A'_w$ . In other-words, by associating an  $O(n^{2\epsilon})$  space structure with every node in every  $\text{GST}(\cdot)$ , document selection queries can be answered in constant time. By replacing  $2\epsilon$  by  $\epsilon$ , we obtain Theorem 7.  $\square$

## 8. Concluding remarks

We have addressed for the first time the problem of finding the  $k$ th document most relevant for a pattern query  $P$  in a collection  $\mathcal{D}$  of string documents. Given the locus node of  $P$  in the generalized suffix tree of  $\mathcal{D}$ , we show how to solve the problem in  $O(\log k)$  time using a linear-space data structure. By using slightly super-linear space,  $O(n \log^\epsilon n)$  for any constant  $\epsilon > 0$ , the time decreases to  $O(1 + \log k / \log \log n)$ . We remark that our structures can be constructed in  $O(n \text{polylog } n)$  time.<sup>1</sup> We also prove that it is highly unlikely that a structure using  $o(n)$  words with query time  $O(\text{polylog } n)$  exists for this problem, and that selection time  $O(1 + \log k / \log \log n)$  is optimal within space  $O(|\text{GST}| \text{polylog } n)$ . The ideas behind our data structure are novel and might find other applications.

A remaining open issue is whether the query time  $O(\log k / \log \log n)$  is optimal in the broader scenario where we can use  $O(n \text{polylog } n)$  space. While for most practical text collections the lower bound is useful as is (since  $|\text{GST}| = \Theta(n)$  in most text families), we believe that it might be significantly harder to obtain the stronger lower bound. It is not even clear if it actually holds, as for example in the case of top- $k$  most frequent elements in a range, the problem on lexicographical

<sup>1</sup> Here we assume that  $w(P, d)$  can be computed in time  $O(t \text{polylog } t)$  time, where  $t$  is the number of occurrences of  $P$  in  $T_d$ .

intervals in a suffix array (which form a hierarchy) can be solved efficiently [32], whereas the problem on arbitrary arrays is most likely harder [3].

### Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

### References

- [1] S. Biswas, A. Ganguly, R. Shah, S.V. Thankachan, Ranked document retrieval for multiple patterns, *Theor. Comput. Sci.* 746 (2018) 98–111.
- [2] S. Biswas, T.-H. Ku, R. Shah, S.V. Thankachan, Position-restricted substring searching over small alphabets, *J. Discret. Algorithms* 46–47 (2017) 36–39.
- [3] T.M. Chan, S. Durocher, K. Green Larsen, J. Morrison, B.T. Wilkinson, Linear-space data structures for range mode query in arrays, *Theory Comput. Syst.* 55 (4) (2014) 719–741.
- [4] Timothy M. Chan, Kasper Green Larsen, Mihai Patrascu, Orthogonal range searching on the ram, revisited, in: *Proceedings of the 27th ACM Symposium on Computational Geometry*, Paris, France, June 13–15, 2011, pp. 1–10.
- [5] Timothy M. Chan, Bryan T. Wilkinson, Adaptive and approximate orthogonal range counting, *ACM Trans. Algorithms* 12 (4) (2016) 45:1–45:15.
- [6] M. Crochemore, C.S. Iliopoulos, M. Kubica, M. Sohel Rahman, G. Tischler, T. Walen, Improved algorithms for the range next value problem and applications, *Theor. Comput. Sci.* 434 (2012) 23–34.
- [7] M.L. Fredman, D.E. Willard, Surpassing the information theoretic barrier with fusion trees, *J. Comput. Syst. Sci.* 47 (1993) 424–436.
- [8] R. Grossi, A. Orlandi, R. Raman, S.S. Rao, More haste, less waste: lowering the redundancy in fully indexable dictionaries, in: *Proc. 26th International Symposium on Theoretical Aspects of Computer Science (STACS)*, 2009, pp. 517–528.
- [9] W.-K. Hon, M. Patil, R. Shah, S.V. Thankachan, J.S. Vitter, Indexes for document retrieval with relevance, in: *Space-Efficient Data Structures, Streams, and Algorithms*, 2013, pp. 351–362.
- [10] W.-K. Hon, M. Patil, R. Shah, S.-B. Wu, Efficient index for retrieving top- $k$  most frequent documents, *J. Discret. Algorithms* 8 (4) (2010) 402–417.
- [11] W.-K. Hon, R. Shah, S.V. Thankachan, J.S. Vitter, On position restricted substring searching in succinct space, *J. Discret. Algorithms* 17 (2012) 109–114.
- [12] W.-K. Hon, R. Shah, S.V. Thankachan, J.S. Vitter, Space-efficient frameworks for top- $k$  string retrieval, *J. ACM* 61 (2) (2014) 9.
- [13] W.-K. Hon, R. Shah, J.S. Vitter, Space-efficient framework for top- $k$  string retrieval problems, in: *Proc. 50th IEEE Symposium on Foundations of Computer Science (FOCS)*, 2009, pp. 713–722.
- [14] Wing-Kai Hon, Sharma V. Thankachan, Rahul Shah, Jeffrey Scott Vitter, Faster compressed top- $k$  document retrieval, in: *2013 Data Compression Conference, DCC 2013, Snowbird, UT, USA, March 20–22, 2013*, pp. 341–350.
- [15] A.G. Jørgensen, K.G. Larsen, Range selection and median: tight cell probe lower bounds and adaptive data structures, in: *Proc. 22nd Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, 2011, pp. 805–813.
- [16] D.T. Lee, C.K. Wong, Worst-case analysis for region and partial region searches in multidimensional binary search trees and balanced quad trees, *Acta Inform.* 9 (1977) 23–29.
- [17] V. Mäkinen, G. Navarro, Position-restricted substring searching, in: *Proc. 7th Latin American Symposium on Theoretical Informatics (LATIN)*, 2006, pp. 703–714.
- [18] J.I. Munro, Tables, in: *Proc. 16th Conference on Foundations of Software Technology and Theoretical Computer Science (FSTTCS)*, 1996, pp. 37–42.
- [19] J.I. Munro, G. Navarro, J.S. Nielsen, R. Shah, S.V. Thankachan, Top- $k$  term-proximity in succinct space, *Algorithmica* 78 (2) (2017) 379–393.
- [20] J.I. Munro, G. Navarro, R. Shah, S.V. Thankachan, Ranked document selection, in: *Proc. 15th Scandinavian Symposium on Algorithmic Theory (SWAT)*, 2014, pp. 344–356.
- [21] J. Ian Munro, Gonzalo Navarro, Jesper Sindahl Nielsen, Rahul Shah, Sharma V. Thankachan, Top- $k$  term-proximity in succinct space, in: *Algorithms and Computation - 25th International Symposium, Proceedings, ISAAC 2014, Jeonju, Korea, December 15–17, 2014*, pp. 169–180.
- [22] S. Muthukrishnan, Efficient algorithms for document retrieval problems, in: *Proc. 13th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, 2002, pp. 657–666.
- [23] G. Navarro, Spaces, trees and colors: the algorithmic landscape of document retrieval on sequences, *ACM Comput. Surv.* 46 (4) (2014), article 52.
- [24] G. Navarro, Y. Nekrich, Top- $k$  document retrieval in optimal time and linear space, in: *Proc. 23rd Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, 2012, pp. 1066–1078.
- [25] G. Navarro, Y. Nekrich, Time-optimal top- $k$  document retrieval, *SIAM J. Comput.* 46 (1) (2017) 89–113.
- [26] G. Navarro, K. Sadakane, Fully-functional static and dynamic succinct trees, *ACM Trans. Algorithms* 10 (3) (2014), article 16.
- [27] G. Navarro, S.V. Thankachan, Faster top- $k$  document retrieval in optimal space, in: *Proc. 20th International Symposium on String Processing and Information Retrieval (SPIRE)*, 2013, pp. 255–262.
- [28] G. Navarro, S.V. Thankachan, New space/time tradeoffs for top- $k$  document retrieval on sequences, *Theor. Comput. Sci.* 542 (2014) 83–97.
- [29] Y. Nekrich, G. Navarro, Sorted range reporting, in: *Proc. 13th Scandinavian Symposium on Algorithmic Theory (SWAT)*, 2012, pp. 271–282.
- [30] Y. Okajima, K. Maruyama, Faster linear-space orthogonal range searching in arbitrary dimensions, in: *Proc. 17th Workshop on Algorithm Engineering and Experiments (ALENEX)*, 2015, pp. 82–93.
- [31] L. Russo, G. Navarro, A. Oliveira, Fully-compressed suffix trees, *ACM Trans. Algorithms* 7 (4) (2011), art. 53.
- [32] R. Shah, C. Sheng, S.V. Thankachan, J.S. Vitter, Top- $k$  document retrieval in external memory, in: *Proc. 21st Annual European Symposium on Algorithms (ESA)*, 2013, pp. 803–814.
- [33] D. Tsur, Top- $k$  document retrieval in optimal space, *Inf. Process. Lett.* 113 (12) (2013) 440–443.
- [34] P. Weiner, Linear pattern matching algorithms, in: *Proc. 14th Annual IEEE Symposium on Switching and Automata Theory*, 1973, pp. 1–11.