



# Suffix-Prefix Queries on a Dictionary



Grigorios Loukides  

Department of Informatics, King's College London, London, UK



Solon P. Pissis  

CWI, Amsterdam, The Netherlands

Vrije Universiteit, Amsterdam, The Netherlands

Sharma V. Thankachan  

North Carolina State University, Raleigh, USA

Wiktor Zuba  

CWI, Amsterdam, The Netherlands

---

## Abstract

---

In the *all-pairs suffix-prefix* (APSP) problem, we are given a dictionary  $R$  of  $k$  strings,  $S_1, \dots, S_k$ , of total length  $n$ , and we are asked to find the length  $\text{SPL}_{i,j}$  of the *longest* string that is both a suffix of  $S_i$  and a prefix of  $S_j$ , for all  $i, j \in [1, k]$ . APSP is a classic problem in string algorithms with many applications in bioinformatics. When all strings of the dictionary are over an integer alphabet of size  $\sigma \leq n^{\mathcal{O}(1)}$ , APSP can be solved in the optimal  $\mathcal{O}(n + k^2)$  time with the use of the generalized suffix tree of the dictionary [Gusfield et al., *Inf. Process. Lett.* 1992].

In many bioinformatics applications, such as in sequence assembly, the size  $k$  of dictionary  $R$  is very large. In particular,  $k^2$  usually dominates  $n$ , and thus the  $k^2$  factor is the bottleneck both in the time and in the space complexity of such applications. We thus initiate a holistic study on several data structure variants of APSP. In particular, we consider the following types of queries:

- **One-to-One**( $i, j$ ): output  $\text{SPL}_{i,j}$ .
- **One-to-All**( $i$ ): output  $\text{SPL}_{i,j}$  for every  $j \in [1, k]$ .
- **Report**( $i, \ell$ ): output all distinct  $j \in [1, k]$  such that  $\text{SPL}_{i,j} \geq \ell$ , where  $\ell \geq 0$  is an integer.
- **Count**( $i, \ell$ ): output the number of distinct  $j \in [1, k]$  such that  $\text{SPL}_{i,j} \geq \ell$ , where  $\ell \geq 0$  is an integer.
- **Top**( $i, K$ ): output  $K$  distinct  $j \in [1, k]$  with the highest values of  $\text{SPL}_{i,j}$  breaking ties arbitrarily.

We assume the standard word RAM model of computation with word size  $w = \Omega(\log n)$  and an integer alphabet of size  $\sigma \leq n^{\mathcal{O}(1)}$ . We show the following upper bounds:

Query	Space (words)	Query time	Note
One-to-One( $i, j$ )	$\mathcal{O}(n)$	$\mathcal{O}(\log \log k)$	Theorem 11
One-to-All( $i$ )	$\mathcal{O}(n)$	$\mathcal{O}(k)$	Theorem 14
Report( $i, \ell$ )	$\mathcal{O}(n)$	$\mathcal{O}(\log n / \log \log n + \text{output})$	Theorem 19(i)
Count( $i, \ell$ )	$\mathcal{O}(n)$	$\mathcal{O}(\log n / \log \log n)$	Theorem 19(ii)
Top( $i, K$ )	$\mathcal{O}(n)$	$\mathcal{O}(\log^2 n / \log \log n + K)$	Theorem 22

We also present efficient algorithms for constructing these data structures.

**2012 ACM Subject Classification** Theory of computation  $\rightarrow$  Pattern matching

**Keywords and phrases** all-pairs suffix-prefix, suffix-prefix queries, internal pattern matching

**Funding** This work is supported in part by the Royal Society grant IES\R3\193209.

*Solon P. Pissis*: Supported in part by the PANGAIA and ALPACA projects that have received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreements No 872539 and 956229, respectively.

*Sharma V. Thankachan*: Supported by the U.S. National Science Foundation (NSF) grant CCF-2146003.

*Wiktor Zuba*: Supported by the Netherlands Organisation for Scientific Research (NWO) through Gravitation-grant NETWORKS-024.002.003.

## 1 Introduction

The *all-pairs suffix-prefix* problem (APSP, in short) is a classic problem in string algorithms. APSP finds numerous applications in bioinformatics because it is the first step in sequence assembly [26, 37, 46, 8, 11]. Given a dictionary  $R$  of  $k$  strings,  $S_1, \dots, S_k$ , of total length  $n$ , the APSP problem asks us to find, for each string  $S_i$ ,  $i \in [1, k]$ , its longest suffix that is a prefix of string  $S_j$ , for all  $j \neq i$ ,  $j \in [1, k]$ . Gusfield et al. [27] presented an algorithm running in the optimal  $\mathcal{O}(n + k^2)$  time for solving APSP, assuming all strings in  $R$  are over an integer alphabet of size  $\sigma \leq n^{\mathcal{O}(1)}$ . The algorithm is based on the generalized suffix tree [53] of  $R$ . Ohlebusch and Gog [39] gave another optimal algorithm which is based on the generalized suffix array [36] of  $R$ . Tustumi et al. [49] gave yet another optimal algorithm based on the generalized suffix array of  $R$ . Thus the common denominator of all existing optimal algorithms for APSP is that they rely on sorting the suffixes of all strings in  $R$ , and therefore they require space  $\Omega(n)$  in any case and for any alphabet. In a very recent work, Loukides and Pissis [34] presented a different optimal algorithm, which is based on the Aho-Corasick automaton of  $R$  [1], and it thus requires space linear in the size of the automaton.

Due to the practical relevance of APSP, there also exists a large body of works devoted to implementing algorithms for APSP that are suboptimal but practically fast on real-world datasets; see [25, 42, 33] and references therein for some of the state-of-the-art implementations. For a parallel implementation of the algorithm by Tustumi et al. see [35]. For approximate variants of APSP, under the Hamming or edit distance, see [44, 52, 32, 5, 47].

In many bioinformatics applications, such as in sequence assembly, the size  $k$  of dictionary  $R$  is very large. In particular,  $k^2$  usually dominates  $n$ , and thus the  $k^2$  factor is the bottleneck both in the time and the space complexity of such applications. For instance, in typical benchmark datasets<sup>1</sup> for genome assembly using short DNA reads (fragments),  $k$  is in the order of  $10^6$  to  $10^8$  and  $n$  is in the order of  $10^8$  to  $10^{10}$ . Hence  $k^2$  dominates  $n$  significantly.

We thus initiate a holistic study on several data structure variants of APSP. Let  $\text{SPL}_{i,j}$  (short for suffix-prefix length), for any  $i, j \in [1, k]$ , denote the length of the *longest* string that is both a suffix of  $S_i$  and a prefix of  $S_j$ . We consider the following types of queries:

- **One-to-One**( $i, j$ ): output  $\text{SPL}_{i,j}$ .
- **One-to-All**( $i$ ): output  $\text{SPL}_{i,j}$  for every  $j \in [1, k]$ .
- **Report**( $i, \ell$ ): output all distinct  $j \in [1, k]$  such that  $\text{SPL}_{i,j} \geq \ell$ , where  $\ell \geq 0$  is an integer.
- **Count**( $i, \ell$ ): output the number of distinct  $j \in [1, k]$  such that  $\text{SPL}_{i,j} \geq \ell$ , where  $\ell \geq 0$  is an integer.
- **Top**( $i, K$ ): output  $K$  distinct  $j \in [1, k]$  with the highest values of  $\text{SPL}_{i,j}$  breaking ties arbitrarily.

By being able to answer different types of such queries efficiently, one may be able to design alternative algorithms, depending on the application in scope, which avoid the  $k^2$  factor in their time or space complexity. Indeed, we stress that most works studying APSP from a practical perspective (e.g., [25, 42, 33]), in fact considered the  $\ell$ -APSP problem in their experimental part; namely, the problem in which we are asked to output only the  $\text{SPL}_{i,j}$  values with  $\text{SPL}_{i,j} \geq \ell$ , for some integer  $\ell \geq 0$ , which, however, is given *a priori* and is *fixed for all pairs*  $S_i, S_j$ . This inflexibility would be surpassed should one have *space-efficient* (e.g., linear-space) data structures for answering these different types of queries *fast*.

<sup>1</sup> For example, see <http://gage.cbcb.umd.edu/data/index.html>.

85 **Our Results** We assume the standard word RAM model of computation with word size  
 86  $w = \Omega(\log n)$  and an integer alphabet of size  $\sigma \leq n^{\mathcal{O}(1)}$ . We show the following upper  
 87 bounds:

Query	Space (words)	Query time	Note
One-to-One( $i, j$ )	$\mathcal{O}(n)$	$\mathcal{O}(\log \log k)$	Theorem 11
One-to-All( $i$ )	$\mathcal{O}(n)$	$\mathcal{O}(k)$	Theorem 14
88 Report( $i, \ell$ )	$\mathcal{O}(n)$	$\mathcal{O}(\log n / \log \log n + \text{output})$	Theorem 19(i)
Count( $i, \ell$ )	$\mathcal{O}(n)$	$\mathcal{O}(\log n / \log \log n)$	Theorem 19(ii)
Top( $i, K$ )	$\mathcal{O}(n)$	$\mathcal{O}(\log^2 n / \log \log n + K)$	Theorem 22

89 We also provide efficient construction algorithms for Theorems 11 and 14: Theorem 11  
 90 can be implemented in  $\mathcal{O}(n \log \log k)$  time and Theorem 14 can be implemented in  $\mathcal{O}(n)$   
 91 time. For Theorems 19 and 22, no guaranteed construction time is provided: the query  
 92 times for Report, Count, and Top rely on the construction of a 2D rectangle stabbing data  
 93 structure for reporting [45] and counting [28], but unfortunately the construction times for  
 94 these data structures are not mentioned in [45] or [28]. However, by constructing the classic  
 95 data structure for 2D rectangle stabbing [15], we obtain  $\mathcal{O}(n \log n)$  construction time,  $\mathcal{O}(n)$   
 96 words of space,  $\mathcal{O}(\log n + \text{output})$  query time for Report,  $\mathcal{O}(\log n)$  query time for Count, and  
 97  $\mathcal{O}(\log^2 n + K)$  query time for Top. We also make the following straightforward observation.

98 **► Observation 1.** *The symmetric versions of One-to-All, Report, Count and Top, where we*  
 99 *are given string  $S_j$  as the query and we are asked to output information about  $\text{SPL}_{i,j}$ , for*  
 100 *all  $i \in [1, k]$ , can be addressed by constructing the corresponding data structures for the*  
 101 *dictionary  $R^r$  of  $k$  strings  $S_1^r, \dots, S_k^r$ , where  $S^r = S[|S|] \cdots S[2]S[1]$  denotes the reverse of*  
 102 *string  $S = S[1]S[2] \cdots S[|S|]$ . Hence, the same space/query-time trade-offs can be achieved.*

103 **Related Work** In addition to the data structure variants of APSP that are studied here,  
 104 two other versions of APSP have been studied in the literature. The first version consists in  
 105 enumerating all pairwise suffix-prefix matches (not necessarily the longest ones) in decreasing  
 106 order of their lengths. This version of the problem was solved by Ukkonen [50], who used  
 107 this solution as the crux of his classic linear-time implementation of the greedy algorithm for  
 108 constructing approximate shortest common superstrings. The second APSP version studied  
 109 consists in enumerating the *set* of longest suffix-prefix matches (not however their association  
 110 with the corresponding pairs of strings) [12]. Since any suffix-prefix match in this set is a  
 111 prefix of some input string, the size of this set is  $\mathcal{O}(n)$ . This version of the problem was  
 112 solved in the optimal  $\mathcal{O}(n)$  time, independently, by Park et al. [40] and by Khan [29].

113 Although our work is inspired by real-world applications, the underlying data structure  
 114 problems are also appealing from a theoretical perspective: (i) they are analogous to *distance*  
 115 *oracles* for networks [48, 41, 17, 16, 13]; and (ii) they are special types of *internal pattern*  
 116 *matching* (IPM) data structures [31, 30, 3, 14, 4]. For instance, an existing, more general,  
 117 IPM data structure [30, 31] can be employed to answer One-to-One queries in  $\mathcal{O}(\log n)$  time  
 118 using  $\mathcal{O}(n)$  words of space; see Section 2.3 for more details. By designing a specialized data  
 119 structure for One-to-One, we obtain  $\mathcal{O}(\log \log k)$  query time using  $\mathcal{O}(n)$  words of space.

120 **Paper Organization** In Section 2, we provide basic definitions and notation on strings. We  
 121 also describe basic data structures for representing a dictionary, some more advanced data  
 122 structures that are necessary to obtain our upper bounds, and a few previous solutions to  
 123 APSP (variants). In Section 3, we provide the solution to One-to-One queries. In Section 4,  
 124 we provide the solution to One-to-All queries. In Section 5, we provide the solutions to Report  
 125 and Count queries. Finally, in Section 6, we provide the solution to Top queries.

## 2 Preliminaries

127 An *alphabet*  $\Sigma$  is a finite nonempty set of  $\sigma = |\Sigma|$  elements called *letters*. By  $\Sigma^*$  we denote  
 128 the set of all strings over  $\Sigma$  including the *empty string*  $\varepsilon$  of length 0. A *string*  $S$  over  $\Sigma$  is a  
 129 sequence of letters of  $\Sigma$ . For a string  $S = S[1] \cdots S[n]$  over  $\Sigma$ , by  $n = |S|$  we denote its length.  
 130 The fragment  $S[i..j]$  of  $S$  is an *occurrence* of the underlying *substring*  $P = S[i] \cdots S[j]$ . We  
 131 also say that  $P$  *occurs* at (*starting*) *position*  $i$  in  $S$ . A *prefix* of  $S$  is a fragment of  $S$  of the  
 132 form  $S[1..j]$  and a *suffix* of  $S$  is a fragment of  $S$  of the form  $S[i..n]$ .

133 Let  $M$  be a finite nonempty set of strings over  $\Sigma$  of total length  $m$ . We call  $M$  a  
 134 *dictionary*. We define the *trie* of  $M$ , denoted by  $\text{TR}(M)$ , as a deterministic finite automaton  
 135 that recognizes  $M$ . Its set of states (nodes) is the set of prefixes of the elements of  $M$ ; the  
 136 initial state (root node) is  $\varepsilon$ ; the set of terminal states is  $M$ ; and transitions (edges) are of the  
 137 form  $\delta(u, \alpha) = u\alpha$ , where  $u$  and  $u\alpha$  are nodes and  $\alpha \in \Sigma$ . The size of  $\text{TR}(M)$  is thus  $\mathcal{O}(m)$ .  
 138 The *compacted trie* of  $M$ , denoted by  $\text{CT}(M)$ , contains the root, the branching nodes, and  
 139 the terminal nodes of  $\text{TR}(M)$ . The term *compacted* refers to the fact that  $\text{CT}(M)$  reduces  
 140 the number of nodes by replacing each maximal branchless path segment with a single edge,  
 141 and that it uses a fragment of a string from  $M$  to represent the label of this edge in  $\mathcal{O}(1)$   
 142 words of space. The nodes of  $\text{TR}(M)$  that are included in  $\text{CT}(M)$  are called *explicit*; all other  
 143 nodes are called *implicit*. The size of  $\text{CT}(M)$  is thus  $\mathcal{O}(|M|)$ . The most well-known form of  
 144 compacted trie is the suffix tree described next.

### 2.1 Suffix Tree and Aho-Corasick Automaton

145 We are given a dictionary  $R$  of  $k$  strings,  $S_1, S_2, \dots, S_k$ , whose total length is  $n = |S_1| + |S_2| +$   
 146  $\cdots + |S_k|$ . Every string in  $R$  is over an integer alphabet  $\Sigma$  whose size  $\sigma$  is polynomial in  $n$ ,  
 147 i.e.,  $\Sigma = \{1, 2, \dots, n^{\mathcal{O}(1)}\}$  and thus  $\sigma \leq n^{\mathcal{O}(1)}$ . For constructing specialized data structures  
 148 and answering internal pattern matching queries, non-trivial representations of  $R$  (different  
 149 than a simple set of strings) are usually more efficient.

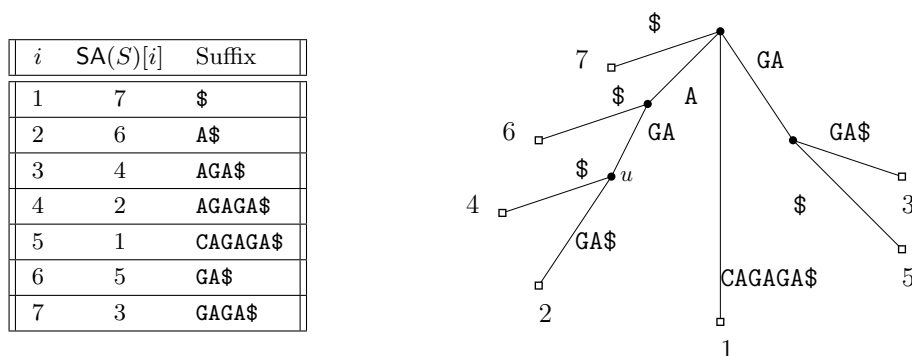
151 Let us set  $T_R := S_1\$1S_2\$2 \cdots S_k\$k$ , where  $\$1 < \$2 < \cdots < \$k$  are letters that are strictly  
 152 lexicographically smaller than any letter from  $\Sigma$  (and as such they do not belong to  $\Sigma$ ).

153 Let  $\text{ST}(S)$  denote the *suffix tree* of string  $S$ , that is the compacted trie of all the suffixes  
 154 of  $S$ . For any node  $v$  of  $\text{ST}(S)$ , by  $\text{str}(v)$  we denote the concatenation of the edge labels on  
 155 the path from the root to  $v$ , and by  $d(v) = |\text{str}(v)|$  we denote the *string depth* of  $v$ . The *suffix*  
 156 *array*  $\text{SA}(S)$  of  $S$  is the lexicographically sorted array of the set of suffixes of  $S$ , represented  
 157 by their starting positions; see Figure 1 for an example.

158 ► **Lemma 2** ([53, 22]). *For any string  $S$  of length  $m$  over an integer alphabet of size*  
 159  *$\sigma \leq m^{\mathcal{O}(1)}$ , the suffix tree and the suffix array of  $S$  can be constructed in  $\mathcal{O}(m)$  time.*

160 We also denote  $\text{ST}_i = \text{ST}(S_i\$i)$  and  $\text{ST}_R = \text{ST}(S_1\$1, \dots, S_k\$k)$ ; that is  $\text{ST}_R$  is the  
 161 generalized suffix tree [51] of the  $k$  strings from  $R$ . The generalized suffix tree can be built  
 162 in linear time; here, however, this more complicated construction is not needed since this  
 163 compacted trie is equivalent to  $\text{ST}(T_R)$  as the letters  $\$i$  occur uniquely in this string (and  
 164 hence a compacted edge containing any label  $\$i$  must end at a leaf node).

165 Another useful representation of  $R$  is given by its Aho-Corasick (AC) automaton [1];  
 166 the set of states of the AC automaton of  $R$ , denoted by  $\text{AC}(R)$ , corresponds to the set  
 167 of the prefixes of the strings in  $R$ . Let  $\text{node}(S)$  denote the node corresponding to string  
 168  $S$ . After reading an input string the automaton must be in a state corresponding to a  
 169 suffix of this string (the longest one that is also a prefix of some string in  $R$  and has  
 170 a corresponding state); such a state always exists as  $\varepsilon$  is always represented (recall  $\varepsilon$  is



■ **Figure 1** Suffix array  $SA(S)$  and suffix tree  $ST(S)$  of string  $S = CAGAGAG\$$ , where  $\$$  is a terminal letter, which is the lexicographically smallest letter occurring in  $S$ . For node  $u$  in  $ST(S)$ ,  $str(u) = AGA$  and  $d(u) = 3$ .

171 the string of length 0). As such, the automaton  $AC(R)$  is often represented by the trie  
 172  $TR(R)$  with transitions  $\delta(\text{node}(S), \alpha) = \{\text{node}(S\alpha)\}$  if  $S\alpha$  is a prefix of a string in  $R$ , and  
 173  $\delta(\text{node}(S), \varepsilon) = \{\text{node}(S')\}$ , where  $S'$  is the longest suffix of  $S$  which is also a prefix of a string  
 174 in  $R$ . The  $\varepsilon$ -transitions are called *failure transitions*. The existence of  $\varepsilon$ -transitions makes  
 175 the automaton nondeterministic, and even though this nondeterminism can be avoided, we  
 176 are going to actually employ those  $\varepsilon$ -transitions to construct the data structure for One-to-All  
 177 queries.

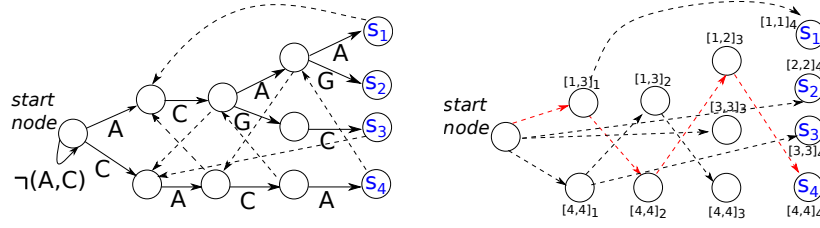
178 ► **Lemma 3** ([1, 20]). *For any dictionary  $R$  of  $k$  strings of total length  $n$  over an integer*  
 179 *alphabet of size  $\sigma \leq n^{\mathcal{O}(1)}$ ,  $AC(R)$  can be constructed in  $\mathcal{O}(n)$  time.*

180 By  $FT(R)$  we denote the so-called *Failure Transition tree* (FTtree) of  $R$ , introduced by  
 181 Loukides and Pissis in [34] for solving the APSP problem: the FTtree nodes correspond to the  
 182 states of the AC automaton (that is, to prefixes of strings in  $R$ ), and the edges correspond to  
 183 its  $\varepsilon$ -transitions with *reversed* direction. Notice that, since every state of  $AC(R)$  has *exactly*  
 184 *one* outgoing failure transition,  $FT(R)$  is indeed a tree rooted at  $\text{node}(\varepsilon)$ . We additionally  
 185 decorate every node  $u$  of  $FT(R)$  by a labeled interval  $I_u = [i, j]_d$ :  $S_i, S_{i+1}, \dots, S_j$  have as  
 186 a common prefix the string of length  $d$  represented by node  $u$ ; see [34]. We will generally  
 187 assume that  $R$  is given lexicographically sorted at construction time; otherwise, the sorted  
 188 version of  $R$  can be produced in linear time using, for example, Lemma 3 or Lemma 2.

189 ► **Example 4.** Let  $R = \{S_1, S_2, S_3, S_4\} = \{ACAA, ACAG, ACGC, CACA\}$  be a dictionary of  $k = 4$   
 190 strings. The AC automaton and the FTtree of  $R$  is shown in Figure 2. Consider the path  
 191 from the root to leaf node  $S_4$  (shown in red) in the FTtree of  $R$ , where the non-root nodes  
 192 have the following labeled intervals  $[i, j]_d$ :  $[1, 3]_1$ ,  $[4, 4]_2$ ,  $[1, 2]_3$ ,  $[4, 4]_4$ . By recording the  
 193 largest string depth  $d$  of an interval containing  $j$ , for every  $j \in [1, k]$ , along this path, we  
 194 compute all  $SPL_{4,j}$ :  $SPL_{4,1} = 3$ ,  $SPL_{4,2} = 3$ ,  $SPL_{4,3} = 1$ , and  $SPL_{4,4} = 4$ . Loukides and  
 195 Pissis [34] showed how to compute this information, for all  $i$ , in  $\mathcal{O}(n + k^2)$  total time, thus  
 196 solving the APSP problem optimally using only the FTtree of  $R$ .

## 197 2.2 Advanced Data Structures

198 Let  $T$  be a rooted tree. A *lowest common ancestor* (LCA) query on  $T$  for two given nodes  
 199  $u$  and  $v$ , denoted by  $w = LCA_T(u, v)$ , returns the last (i.e., the lowest) common node  $w$  on  
 200 their paths from the root.



■ **Figure 2** The AC automaton  $AC(R)$  (on the left) and FTtree  $FT(R)$  (on the right) of the dictionary of strings  $R = \{S_1, S_2, S_3, S_4\} = \{ACAA, ACAG, ACGC, CACA\}$ . In  $AC(R)$ , solid arrows correspond to transitions and dashed arrows to failure transitions. To avoid cluttering the figure, failure transitions to the start node in  $AC(R)$  have been omitted.

201 ► **Lemma 5** ([9]). *For any rooted tree  $T$  with  $m$  nodes, after  $\mathcal{O}(m)$ -time preprocessing, we*  
 202 *can answer  $LCA_T$  queries in  $\mathcal{O}(1)$  time per query.*

203 A *rank and select* data structure (also known as *succinct indexable dictionary* [43]) is a  
 204 classic data structure, constructed over an array  $A$  of length  $m$  over alphabet  $[1, \sigma]$ , which  
 205 supports two types of queries:

- 206 ■  $\text{rank}_A(i, x) = |\{\ell \in [1, x] : A[\ell] = i\}|$ , for  $i \in [1, \sigma]$  and  $x \in [1, m]$ ;
- 207 ■  $\text{select}_A(i, x) = \min\{\ell \in [1, m] : \text{rank}_A(i, \ell) = x\}$ , for  $i \in [1, \sigma]$  and  $x \in [1, m]$ .

208 In other words,  $\text{rank}_A(i, x)$  returns the number of elements with value equal to  $i$  occurring at  
 209 positions in  $[1, x]$  of  $S$ , while  $\text{select}_A(i, x)$  returns the position of the  $x$ th element of  $A$  with  
 210 value equal to  $i$ .

211 ► **Lemma 6** ([7, 38, 18]). *For any array  $A = A[1..m]$  over  $[1, \sigma]$ ,  $\sigma \leq m$ , after  $\mathcal{O}(m \log \log \sigma)$ -*  
 212 *time preprocessing, we can construct a data structure of  $\mathcal{O}(m)$  words of space that supports*  
 213  *$\mathcal{O}(\log \log \sigma)$ -time rank and select queries on  $A$ .*

214 Let  $T$  be a rooted tree of  $m$  nodes with integer weights on nodes. Further assume that  
 215 the weight of every node of  $T$  satisfies the *min-heap property*: the weight of each node is  
 216 greater than or equal to the value of its parent (the smallest weight is hence at the root).  
 217 A *weighted ancestor* (WA) query for a given node  $u$  of  $T$  and an integer  $d$ , denoted by  
 218  $w = \text{WA}_T(u, d)$ , returns its deepest ancestor  $w$  whose weight is at most  $d$  [23]. This problem  
 219 is the generalization of the classic *predecessor search* problem on rooted trees. In the special  
 220 case when  $T$  is a suffix tree and the nodes are weighted by *string depth*, the problem admits  
 221 an optimal solution due to the recent result of Belazzougui et al. [6] (see also [24]).

222 ► **Lemma 7** ([6]). *For any suffix tree  $T$  with  $m$  nodes weighted by string depth, after*  
 223  *$\mathcal{O}(m)$ -time preprocessing, we can answer  $WA_T$  queries in  $\mathcal{O}(1)$  time per query.*

224 In this special case, the ancestor at string depth exactly  $d$  may be an implicit node of  $T$ ,  
 225 in which case the query outputs its closest explicit ancestor.

## 2.3 Previous Solutions

227  $\mathcal{O}(n + k^2)$ -**time Algorithm for APSP** We describe the optimal solution to APSP given by  
 228 Gusfield et al. in [27]. We set  $T_R := S_1 \$_1 S_2 \$_2 \cdots S_k \$_k$ , where  $\$_1 < \$_2 < \cdots < \$_k$  are letters  
 229 that are strictly lexicographically smaller than any letter from  $\Sigma$ . We start by constructing  
 230 the suffix tree  $ST_R = ST(T_R)$ . Using a DFS traversal on  $ST_R$ , we construct lists  $L(v)$  for  
 231 all nodes  $v$  of  $ST_R$ :  $L(v)$  stores all  $i$  such that the suffix of length  $d(v)$  of string  $S_i$  is  $\text{str}(v)$ .

232 Consider a string  $S_j$  from  $R$  and focus on the path  $P_j$  from the root of  $\text{ST}_R$  to the leaf node  
 233 representing the longest suffix of  $S_j$ , i.e., the entire string  $S_j$ . Let  $v$  be a node on  $P_j$ . A  
 234 suffix of string  $S_i$  of length  $d(v)$  is a prefix of string  $S_j$  of the same length if and only if  $i$  is  
 235 in  $L(v)$ . However, for each index  $i$ , we want to record the *deepest* node  $v$  on  $P_j$  such that  $i$  is  
 236 in  $L(v)$ . It then follows that  $d(v) = \text{SPL}_{i,j}$ . In order to achieve a linear-time complexity, we  
 237 perform another DFS maintaining  $k$  stacks (one for each  $S_i$ ). Upon visiting  $v$ , we push it on  
 238 stack  $i$  for every  $i \in L(v)$ . When the leaf node representing the entire string  $S_j$  is reached,  
 239 we scan the  $k$  stacks and record, for each index  $i$ , the current top of the  $i$ th stack. When  $v$  is  
 240 reached in a backward edge traversal, we pop the top of any stack whose index is in  $L(v)$ .  
 241 We obtain the following result.

242 ► **Lemma 8** ([27]). *For any dictionary of  $k$  strings of total length  $n$  over an integer alphabet*  
 243 *of size  $\sigma \leq n^{\mathcal{O}(1)}$ , APSP can be solved in the optimal  $\mathcal{O}(n + k^2)$  time.*

244 In what follows, we assume that  $k \geq \sqrt{n}$ ; otherwise, when  $k < \sqrt{n}$ , Lemma 8 implies an  
 245 optimal solution to our data structure problems (linear preprocessing time, linear size and  
 246 time-optimal queries), which precomputes and stores all answers.

247 **Internal Prefix-Suffix Queries for One-to-One** Kociumaka considered the following data  
 248 structure problem in [30]: Given two fragments  $x$  and  $y$  of a string  $T$  and a positive integer  
 249  $d$ , report all suffixes of  $y$  of length between  $d$  and  $2d - 1$  that also occur as prefixes of  $x$   
 250 (represented as an arithmetic progression of their lengths). This is the *Internal Prefix-Suffix*  
 251 *Queries* problem. Kociumaka showed the following result (see also [31]).

252 ► **Lemma 9** (Theorem 1.1.3 in [30]). *For any string  $T$  of length  $m$  over an integer alphabet of*  
 253 *size  $\sigma \leq m^{\mathcal{O}(1)}$ , after  $\mathcal{O}(m)$ -time preprocessing, we can answer Internal Prefix-Suffix Queries*  
 254 *in  $\mathcal{O}(1)$  time per query.*

255 By employing Lemma 9 on  $T_R$ , after an  $\mathcal{O}(n)$ -time preprocessing, we can answer  
 256 **One-to-One** queries in  $\mathcal{O}(\log(\min(|S_i|, |S_j|))) = \mathcal{O}(\log n)$  time. In particular, we query  
 257 for  $x = S_j$ ,  $y = S_i$ , and  $d = 2^\ell$ , for all integers  $0 \leq \ell \leq \log \min(|S_i|, |S_j|)$ , to compute a  
 258 representation of all the suffixes of  $S_i$  that are also prefixes of  $S_j$  and then return the length  
 259 of the longest one as  $\text{SPL}_{i,j}$ . We obtain the following result, which we improve in Section 3.

260 ► **Corollary 10.** *For any dictionary of  $k$  strings of total length  $n$  over an integer alphabet*  
 261 *of size  $\sigma \leq n^{\mathcal{O}(1)}$ , we can construct a data structure of  $\mathcal{O}(n)$  words of space answering*  
 262 *One-to-One queries in  $\mathcal{O}(\log n)$  time.*

### 263 **3 Answering One-to-One Queries**

264 **Main Idea** Say we want to find the longest suffix of  $S_i$  that is a prefix of  $S_j$ . We first  
 265 find the maximal longest common prefix between  $S_j$  and any suffix of  $S_i$ . Say this suffix  
 266 is  $S_i[q..|S_i|]$  and we have that  $S_i[q..q+r-1] = S_j[1..r]$  is this longest common prefix.  
 267 If this prefix is the whole  $S_i[q..|S_i|]$ , i.e.,  $|S_i| = q+r-1$ , then  $r$  is clearly the answer. If  
 268 this longest common prefix is not a suffix of  $S_i$ , i.e.,  $|S_i| > q+r-1$ , then the answer is the  
 269 longest prefix of  $S_i[q..q+r-1]$ , that is also a suffix of  $S_i$ .

270 Recall that  $\text{ST}_i = \text{ST}(S_i\$i)$  and  $\text{ST}_R = \text{ST}(T_R)$ . Consider the path in  $\text{ST}_R$  obtained by  
 271 reading  $S_j\$j$  from its root (this path ends in a leaf node). When spelling any suffix of  $S_i$   
 272 that is also a prefix of  $S_j$  in  $\text{ST}_R$  we use exactly the same path and end by going out of it  
 273 when reading  $\$i$ . This means, that  $\text{SPL}_{i,j}$  is represented by the lowest node on this path that  
 274 has an outgoing edge with label  $\$i$ .

275 In the following we focus on enhancing  $\text{ST}_R$  and  $\text{ST}_i$ , for all  $i \in [1, k]$ , to obtain a data  
 276 structure that allows finding the string depth of such a node (equal to  $\text{SPL}_{i,j}$ ) efficiently. We  
 277 will prove the following result.

278 **► Theorem 11.** *For any dictionary of  $k$  strings of total length  $n$  over an integer alphabet of size  
 279  $\sigma \leq n^{\mathcal{O}(1)}$ , we can construct a data structure of  $\mathcal{O}(n)$  words of space answering **One-to-One**  
 280 queries in  $\mathcal{O}(\log \log k)$  time. The data structure can be constructed in  $\mathcal{O}(n \log \log k)$  time.*

281 Let us start with a straightforward auxiliary lemma.

282 **► Lemma 12.** *For any dictionary of  $k$  strings  $S_1, \dots, S_k$  of total length  $n$  over an integer  
 283 alphabet of size  $\sigma \leq n^{\mathcal{O}(1)}$ , in  $\mathcal{O}(n)$  time we can construct a data structure of  $\mathcal{O}(k)$  words of  
 284 space that answers queries of the type “Is  $S_j$  a suffix of  $S_i$ ?” in  $\mathcal{O}(1)$  time.*

285 **Proof.** Let  $X^r$  denote the reverse of string  $X$ , i.e.,  $X^r = X[|X|] \dots X[1]$ . We first sort  
 286  $S_1^r, \dots, S_k^r$  lexicographically, and store for each  $j \in [1, k]$  a value  $\text{rlex}[j] \in [1, k]$  equal to the  
 287 rank of  $S_j^r$  in this sorted list.  $S_j$  is a suffix of  $S_i$  if and only if  $S_j^r$  is a prefix of  $S_i^r$ . The  
 288 crucial property of this ordering is that all the strings such that  $S_j^r$  is their prefix form an  
 289 interval from the position  $\text{rlex}[j]$  to a position  $\text{rlex}[j] + l[j] - 1$ , where  $l[j]$  is the total number  
 290 of strings  $S_1^r, \dots, S_k^r$  starting with  $S_j^r$ ; that is,  $\text{rlex}[j] + l[j]$  is the position of the first string  
 291 having a longest common prefix with  $S_j^r$  shorter than  $|S_j^r|$ . The values  $\text{rlex}[j]$  and  $l[j]$ , for all  
 292  $j \in [1, k]$ , can be computed in  $\mathcal{O}(n)$  time [19].

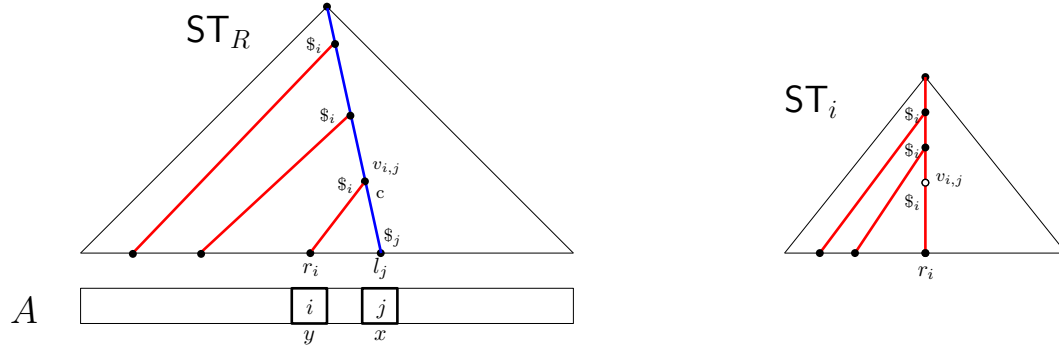
293 As for the querying, for any  $i, j$ , we have that  $S_j$  is a suffix of  $S_i$  if and only if  $\text{rlex}[j] \leq$   
 294  $\text{rlex}[i] < \text{rlex}[j] + l[j]$ , which is checked in  $\mathcal{O}(1)$  time. The total size of arrays  $l$  and  $\text{rlex}$  is  
 295  $\Theta(k)$ . ◀

296 **Construction** We start the construction of the data structure by constructing the data struc-  
 297 ture underlying Lemma 12. We also construct  $\text{ST}_R$  and  $\text{ST}_i$ , for all  $i \in [1, k]$ , using Lemma 2.  
 298 We enhance  $\text{ST}_R$  with the data structure for LCA queries underlying Lemma 5, and link  
 299 the leaf nodes originating from suffixes of  $S_i\$i$  with the corresponding leaf nodes of  $\text{ST}_i$ , for  
 300 all  $i \in [1, k]$ . We construct an array  $A = A[1..|T_R|]$  over  $[1, k]$  such that  $A[\ell] = i$  if the  $\ell$ th  
 301 leaf node (from the left) of  $\text{ST}_R$  originates from a suffix of  $S_i\$i$ ; since the leaf nodes are  
 302 ordered according to the lexicographic order of the suffixes they originate from, array  $A$  can  
 303 be easily extracted from  $\text{SA}(T_R)$  constructed by means of Lemma 2. We enhance array  $A$   
 304 with the rank and select data structure underlying Lemma 6. We link the leaf nodes of  $\text{ST}_R$   
 305 with the corresponding elements of  $A$ . For each  $\text{ST}_i$ , we construct the data structure for  
 306 WA queries underlying Lemma 7. For every node  $w$  of  $\text{ST}_i$ , we store the string depth of  
 307 its closest ancestor (including  $w$  itself) that has an outgoing edge with label  $\$i$  and hence  
 308 corresponds to a suffix of  $S_i$ ; since the root always has such an edge, this assignment is always  
 309 well-defined. In order to efficiently compute and store all those values, we simply process the  
 310 information through the tree in a top-down manner. This completes the construction.

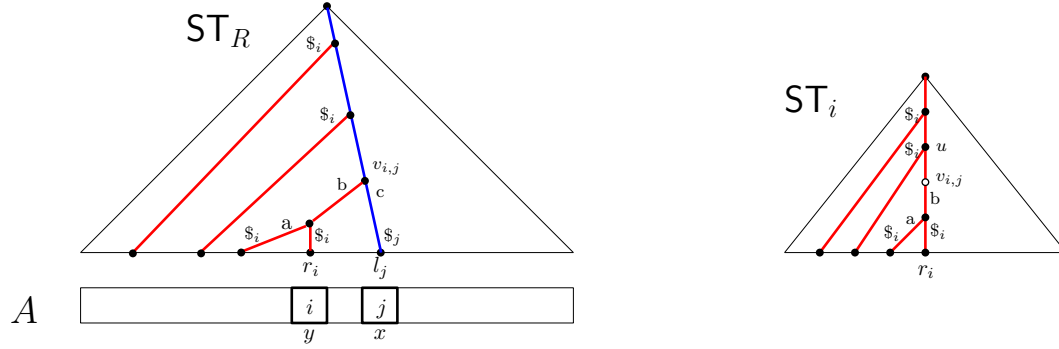
311 The part of the data structure that relies on Lemmas 2, 5, 7, and 12 is implemented in  
 312  $\mathcal{O}(n)$  time and it occupies  $\mathcal{O}(n)$  words of space. By Lemma 6, array  $A$  occupies  $\mathcal{O}(n)$  words  
 313 of space, and it can be implemented in  $\mathcal{O}(n \log \log k)$  time as it stores  $k$  distinct values.

314 **Querying** Consider a **One-to-One**( $i, j$ ) query; that is, we want to compute  $\text{SPL}_{i,j}$ , the length  
 315 of the longest suffix of  $S_i$  that is a prefix of  $S_j$ . Let  $x$  be the position in array  $A$  that  
 316 corresponds to the leaf node  $l_j$  of  $\text{ST}_R$  reached after conceptually reading  $S_j\$j$ . We first  
 317 check if the entire  $S_j$  is a suffix of  $S_i$  by means of Lemma 12. If this is the case then we return  
 318  $\text{SPL}_{i,j} = |S_j|$ . If this is not the case (inspect Figure 3), we perform the following sequence





■ **Figure 3** An illustration of the  $\text{One-to-One}(i, j)$  query algorithm. The node  $v_{i,j}$ , which is explicit in  $\text{ST}_R$  but implicit in  $\text{ST}_i$ , has an outgoing edge labeled with  $\$i$  and hence the string depth  $d(v_{i,j})$  of node  $v_{i,j}$  is the answer to the query.



■ **Figure 4** An illustration of the  $\text{One-to-One}(i, j)$  query algorithm. The closest ancestor of node  $v_{i,j}$ , which is explicit in  $\text{ST}_R$  but implicit in  $\text{ST}_i$ , with an outgoing edge labeled with  $\$i$  is node  $u$  and hence the string depth  $d(u)$  of node  $u$  is the answer to the query.

319 of queries,  $\text{select}_A(i, \text{rank}_A(i, x))$ , which finds the position  $y$  in array  $A$  that corresponds to  
 320 the leaf node  $r_i$ ; this corresponds to the suffix of  $S_i\$i$  that is closest to the left of  $l_j$ . We  
 321 then compute the lowest common ancestor of  $r_i$  and  $l_j$ :  $v_{i,j} = \text{LCA}_{\text{ST}_R}(r_i, l_j)$ . If node  $v_{i,j}$   
 322 has an outgoing edge labeled with  $\$i$ , which ends at  $r_i$ , then we return  $\text{SPL}_{i,j} = d(v_{i,j})$  (this  
 323 is the case in Figure 3). We check this by checking whether  $d(r_i) = d(v_{i,j}) + 1$ . If  $v_{i,j}$  does  
 324 not have such an outgoing edge (this is the case in Figure 4), we locate the explicit node  
 325 corresponding to  $v_{i,j}$  in  $\text{ST}_i$  (or its closest explicit ancestor if it is implicit) by asking a WA  
 326 query:  $w = \text{WA}_{\text{ST}_i}(r_i, d(v_{i,j}))$ . Finally, we return the string depth of the closest ancestor  
 327 of  $w$  with an outgoing edge labeled  $\$i$  as  $\text{SPL}_{i,j}$ ; recall that every node of  $\text{ST}_i$  stores this  
 328 information.

329 The time complexity of the query is  $\mathcal{O}(\log \log k)$ ; the bottleneck is the complexity of the  
 330 rank and select queries on  $A$  – all other operations take constant time. Let us now explain  
 331 why the faster  $\mathcal{O}(1)$ -time select and  $\mathcal{O}(1 + \log \frac{\log k}{\log w})$ -time rank queries presented in [7], where  
 332  $w$  is the machine word, cannot improve our query time further. The size of the problem  
 333 is  $\Theta(n)$ , hence the size of the machine word in the word-RAM model is  $\Theta(\log n)$ , thus the  
 334 query time equals  $\mathcal{O}(1 + \log \frac{\log k}{\log \log n})$ . However, we have assumed that  $k \geq \sqrt{n}$  (otherwise  
 335 the structure of Lemma 8 implies an optimal solution – linear size and constant time queries  
 336 – for the One-to-One queries), hence this is equal to  $\mathcal{O}(1 + \log \log k) = \mathcal{O}(\log \log k)$  as stated.

337 **Correctness** Recall that the answer to `One-to-One( $i, j$ )` equals to the string depth of the  
 338 closest ancestor of  $l_j$  in  $\text{ST}_R$  that has an outgoing edge labeled with  $\$_i$ . By construction,  
 339 this ancestor ends on the right of  $l_j$  only if the entire  $S_j$  is a suffix of  $S_i$ , which we check  
 340 separately. Otherwise, this ancestor is also an ancestor of  $r_i$  (which is on the left of  $l_i$ ) as  $\$_i$   
 341 goes out of the path from the root to  $l_j$  to the left (by construction, it is lexicographically  
 342 smaller than the next letter on this path), and hence this edge labeled with  $\$_i$  must end  
 343 either in  $r_i$  or further to the left (by the definition of  $r_i$ ). As an ancestor of  $l_j$  and  $r_i$ , it is also  
 344 the closest ancestor of  $v_{i,j}$  with such an outgoing edge; the latter actually exists (possibly as  
 345 an implicit node) in  $\text{ST}_i$  (unlike  $l_j$ ). The final steps of the query algorithm find the string  
 346 depth of the node corresponding to the searched ancestor in  $\text{ST}_i$  (string depth is a shared  
 347 property of the corresponding nodes).

348 We have arrived at Theorem 11. Note that the construction time for our data structure  
 349 is  $\mathcal{O}(n \log \log k)$ . The bottleneck for the construction time is the construction time for the  
 350 rank and select data structure (Lemma 6).

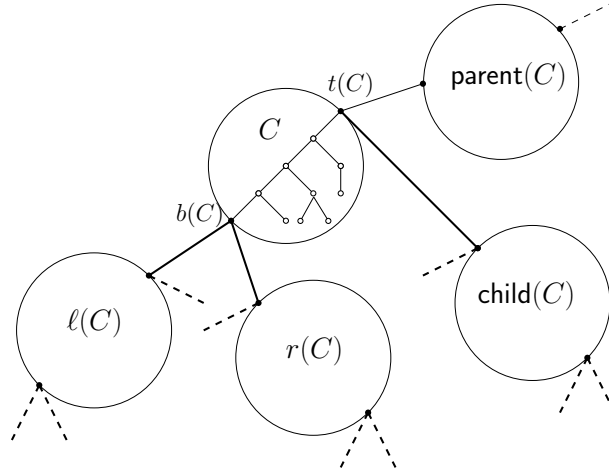
## 351 4 Answering One-to-All Queries

352 The spine of the data structure described in this section is  $\text{FT}(R)$ , the FTtree of  $R$  (see Sec-  
 353 tion 2). Recall that for each node in  $\text{FT}(R)$  (representing each prefix of a string  $S_i$ ), we store  
 354 information about which strings from  $R$  it is a prefix of (see Figure 2).

355 **Main Idea** The Aho-Corasick lemma [1] states that for any two nodes,  $\text{node}(U)$  and  $\text{node}(V)$ ,  
 356 in  $\text{AC}(R)$ , we have a failure transition from  $\text{node}(U)$  to  $\text{node}(V)$  if and only if  $V$  is the longest  
 357 suffix of  $U$  that is also a prefix of some string in  $R$ . As a consequence, in  $\text{FT}(R)$ ,  $\text{node}(S)$   
 358 is an ancestor of  $\text{node}(S')$  if and only if  $S$  is a suffix of  $S'$  (and both are prefixes of some  
 359 strings from  $R$  as nodes of  $\text{FT}(R)$ ). Thus the path from  $\text{node}(\varepsilon)$  (the root) to  $\text{node}(S_i)$  in  
 360  $\text{FT}(R)$  contains exactly the nodes  $\text{node}(S)$  such that  $S$  is a suffix of  $S_i$  and a prefix of some  
 361 string in  $R$ . Those nodes are ordered according to the string length, hence the nodes closer  
 362 to  $\text{node}(S_i)$  on this path will correspond to *longer* suffix-prefix matches.

363 A `One-to-All( $i$ )` query can thus be answered by simply reading the path from the root to  
 364  $\text{node}(S_i)$  recording, for each  $j \in [1, k]$ , the last node on the path corresponding to a prefix of  
 365  $S_j$ . The space occupied by  $\text{FT}(R)$  is in  $\mathcal{O}(n)$ ; and such a query algorithm can take  $\Theta(|S_i|)$ ,  
 366 that is even  $\Theta(n)$  time. Hence, by such an algorithm, we would not really gain anything from  
 367 constructing  $\text{FT}(R)$  in the preprocessing. On the other extreme, by running this algorithm  
 368 not for a single path, but for the whole  $\text{FT}(R)$  using a DFS traversal, we can precompute the  
 369 answers for all the values of  $i \in [1, k]$  in  $\mathcal{O}(n + k^2)$  total time (and space), and then answer a  
 370 query in  $\mathcal{O}(k)$  time by simply outputting the  $k$  stored values; this would not be faster than  
 371 using the algorithm by Gusfield et al. [27] or the one by Loukides and Pissis [34]. We will  
 372 augment  $\text{FT}(R)$  to obtain a more efficient solution combining the space efficiency of the first  
 373 approach with the low query time of the second one.

374 A  $\tau$ -*micro-macro decomposition*, introduced for rooted binary trees in [2], and then  
 375 generalized for rooted general trees in [10] (after an appropriate mapping), is a partition of a  
 376 rooted tree  $T$  of  $N$  nodes into  $\mathcal{O}(N/\tau)$  connected subtrees, called *micro trees*. In the case of  
 377 binary trees each micro tree is of size at most  $\tau$  and at most two of its nodes are adjacent to  
 378 nodes in other micro trees. These nodes are referred to as *top* and *bottom boundary* nodes  
 379 of the micro tree. The top boundary node is chosen as the root of the micro tree. The  
 380 *macro tree* is a rooted tree of size  $\mathcal{O}(N/\tau)$  whose nodes correspond to micro trees as follows  
 381 (inspect Figure 5): The top boundary node  $t(C)$  of a micro tree  $C$  is connected to a boundary



■ **Figure 5** The structure of a micro-macro decomposition of a rooted binary tree.

node  $\text{parent}(C)$  in the parent micro tree (apart from the root). The boundary node  $t(C)$  might also be connected to a top boundary node of a child micro tree, which we denote by  $\text{child}(C)$ . Such a  $\tau$ -micro-macro decomposition can be computed in  $\mathcal{O}(N)$  time for binary [2] and general [10] rooted trees. We summarize the above discussion in the lemma below.

► **Lemma 13** ([2, 10]). *For any rooted tree  $T$  with  $N$  nodes and for any integer  $\tau \in [1, N]$ , the  $\tau$ -micro-macro decomposition of  $T$  can be computed in  $\mathcal{O}(N)$  time.*

We will prove the following result.

► **Theorem 14.** *For any dictionary of  $k$  strings of total length  $n$  over an integer alphabet of size  $\sigma \leq n^{\mathcal{O}(1)}$ , we can construct a data structure of  $\mathcal{O}(n)$  words of space answering One-to-All( $i$ ) queries in  $\mathcal{O}(k)$  time. The data structure can be constructed in  $\mathcal{O}(n)$  time.*

**Construction** We start the construction of the data structure by constructing  $\text{FT}(R)$  from  $\text{AC}(R)$  using Lemma 3. We compute the  $\tau$ -micro-macro decomposition of  $\text{FT}(R)$ , for a parameter  $\tau$  defined later, using Lemma 13. For each node  $u$  of the  $\text{FT}(R)$ , corresponding to a prefix  $S$  of some string  $S_i$  in  $R$ , we store the labeled interval  $I_u$ . For each boundary node in the  $\tau$ -micro-macro decomposition of  $\text{FT}(R)$ , we store an array of  $k$  integers, which for each  $i \in [1, k]$ , stores the string depth of its lowest ancestor  $\text{node}(S)$  such that  $S$  is a prefix of  $S_i$ . The additional size for storing this information in all the boundary nodes is  $\mathcal{O}(k \cdot n/\tau)$ . We compute these arrays by performing a DFS over  $\text{FT}(R)$  with a set of  $k$  stacks, one for every string in  $R$ , storing the string depths of ancestors of the visited node of each type (which  $S_i$  they originate from). As there are only  $2n$  updates of the stacks (each prefix of a string  $S_i$  is stored and removed once from the  $i$ th stack) and the information is stored by simply reading the top values of the  $k$  stacks, the total computation time is bounded by  $\mathcal{O}(n + k \cdot n/\tau)$ .

**Querying** Let us start with the following observation from [34] (inspect also Figure 2).

► **Observation 15** ([34]). *Let  $u$  and  $v$  be two non-root nodes of  $\text{FT}(R)$  with labeled intervals  $I_u = [i_u, j_u]_{d(u)}$  and  $I_v = [i_v, j_v]_{d(v)}$ , respectively, and such that  $u$  is an ancestor of  $v$ . Then  $d(u) < d(v)$  and either  $[i_u, j_u]$  contains  $[i_v, j_v]$  or  $[i_u, j_u]$  and  $[i_v, j_v]$  do not intersect.*

408 Consider a One-to-All( $i$ ) query; that is, we want to compute an array of length  $k$ , which  
 409 stores  $\text{SPL}_{i,j}$ , for all  $j \in [1, k]$ . We start by finding the closest boundary node on the path  
 410 from the root to  $\text{node}(S_i)$ ; that is, the top boundary node of the micro tree containing  
 411  $\text{node}(S_i)$ . On the path between this top boundary node and  $\text{node}(S_i)$ , there are at most  $\tau$   
 412 nodes. We compute the information coming from just those nodes in  $\mathcal{O}(k + \tau)$  time with a  
 413 sweep line approach: there are  $\mathcal{O}(\tau)$  (labeled) intervals from  $[1, k]$ , the intervals are labeled  
 414 by different values (string depth), but, by Observation 15, two intervals are either disjoint  
 415 or the one with the larger string depth is contained in the one with the smaller one. Thus,  
 416 it is enough to hold the active intervals on a stack to keep track of the longest possible  
 417 suffix-prefix match: the interval on the top of the stack has the highest value and will end  
 418 the soonest. The solution is then obtained as the position-wise maximum of the computed  
 419 array and the array stored in the top boundary node, which we compute in  $\mathcal{O}(k)$  time.

420 **Correctness** The correctness of the algorithm follows by the Aho-Corasick lemma (see also  
 421 the discussion of the “main idea” paragraph above).

422 The data structure occupies  $\mathcal{O}(n + k \cdot n/\tau)$  words of space and supports One-to-All queries  
 423 in  $\mathcal{O}(k + \tau)$  time. By setting  $\tau$  to  $k$  (or to  $ck$ , for some positive constant  $c$  that balances the  
 424 operation costs more efficiently) we obtain the complexities claimed in Theorem 14. Note  
 425 that the data structure is constructed in  $\mathcal{O}(n + k \cdot n/\tau)$  time, which is  $\mathcal{O}(n)$  for  $\tau = \Theta(k)$ .  
 426 Thus the presented data structure for One-to-All queries is optimal.

## 427 5 Answering Report and Count Queries

428 In this section we are going to use  $\text{ST}_R$  again. This time, however, instead of augmenting  
 429  $\text{ST}_R$  with an LCA data structure and linking its nodes with the rank and select array, we  
 430 are going to link the nodes with rectangles and employ classic results from computational  
 431 geometry for reporting (see Lemma 16) and counting (see Lemma 17).

432 Let  $[x_1, x_2] \times [y_1, y_2]$  denote a rectangle in a 2D space with edges parallel to the axes,  
 433 where the intervals  $[x_1, x_2]$  and  $[y_1, y_2]$  are the projections of this rectangle to the  $x$ -axis and  
 434  $y$ -axis, respectively. In the reporting version of the *2D rectangle stabbing* problem [15], we are  
 435 given a set  $S$  of  $n$  rectangles to preprocess, so that when we are given a query point  $q = (x, y)$ ,  
 436 we report the subset  $Q \subseteq S$  of rectangles  $[x_1, x_2] \times [y_1, y_2]$  that contain  $q$ :  $x_1 \leq x \leq x_2$  and  
 437  $y_1 \leq y \leq y_2$ . In the counting version of 2D rectangle stabbing, we are asked to return  $|Q|$ .

438 **► Lemma 16** ([45]). *For any set  $S$  of  $n$  rectangles, we can construct a data structure of  $\mathcal{O}(n)$   
 439 words of space answering 2D rectangle stabbing reporting queries in  $\mathcal{O}(\log n / \log \log n + f)$   
 440 time, where  $f$  is the output size  $|Q|$ .*

441 2D rectangle stabbing counting is known to be reducible to 2D orthogonal range count-  
 442 ing [21], and such a data structure for 2D orthogonal range counting can be found in [28].

443 **► Lemma 17** ([21, 28]). *For any set  $S$  of  $n$  rectangles, we can construct a data structure of  
 444  $\mathcal{O}(n)$  words of space answering 2D rectangle stabbing counting queries in  $\mathcal{O}(\log n / \log \log n)$   
 445 time.*

446 **Main Idea** For every suffix  $S$  of a string in  $R$  that is represented by a node in  $\text{ST}_R$ , we  
 447 define a rectangle in 2D space: the  $x$  dimension corresponds to the lexicographically sorted  
 448 list of all suffixes of strings in  $R$  whose prefix is  $S$ ; and the  $y$  dimension corresponds to  
 449 interval  $[0, |S|]$ . A Report (resp. a Count) query is defined by two parameters, which form a  
 450 point in the 2D space:  $i$  corresponds to string  $S_i$  in the same sorted list ( $x$  dimension) and  $\ell$

451 corresponds to the smallest length of interest ( $y$  dimension). By reporting (resp. counting)  
 452 all rectangles *enclosing this point* (Lemmas 16 and 17), we locate all suffix-prefix matches.  
 453 Extra care, however, needs to be taken in order to avoid double reporting (resp. counting).

454 **Construction** We start the construction of the data structure by constructing  $\text{ST}_R$  us-  
 455 ing Lemma 2. Let  $u$  be an explicit or implicit node of  $\text{ST}_R$  that is the parent of a leaf node  
 456 reached with  $\$_i$ : the labels of the path from root to  $u$  form a suffix of  $S_i$ . For every such node  
 457  $u$  and every  $i$ , we create a tuple  $(L(u), R(u), d(u), i)$ , where  $L(u)$  and  $R(u)$  are the (pre-order  
 458 rank of) the leftmost and the rightmost leaf node under  $u$ , respectively.<sup>2</sup> Note that such a  
 459 node may correspond to multiple tuples for different  $i$  values – this occurs when distinct  
 460 elements of  $R$  share the same suffix. There are exactly  $n$  such tuples (one for every suffix)  
 461 coming from  $\text{ST}_R$  and we can compute them in  $\mathcal{O}(n)$  total time using a DFS traversal.

462 Recall that if we spell  $S_j\$_j$  in  $\text{ST}_R$  and the obtained leaf node  $v$  has an ancestor of string  
 463 depth  $\ell$  which has an outgoing edge with label  $\$_i$ , then  $\text{SPL}_{i,j} \geq \ell$ . The same property  
 464 ( $\text{SPL}_{i,j} \geq \ell$ ) can be expressed by  $L(v) \in [L(u), R(u)]$  (namely,  $u$  is an ancestor of  $v$ ), and  
 465  $\ell \in [0, d(u)]$  (namely, the string depth of  $u$  is at least  $\ell$ ) for a tuple  $(L(u), R(u), d(u), i)$ .  
 466 Now note that  $(L(u), R(u), d(u), i)$  forms a rectangle, whose identifier is  $i$ . In particular,  
 467  $(L(u), R(u), d(u), i)$  can be viewed as rectangle  $[L(u), R(u)] \times [0, d(u)]$  with satellite data  $i$ .

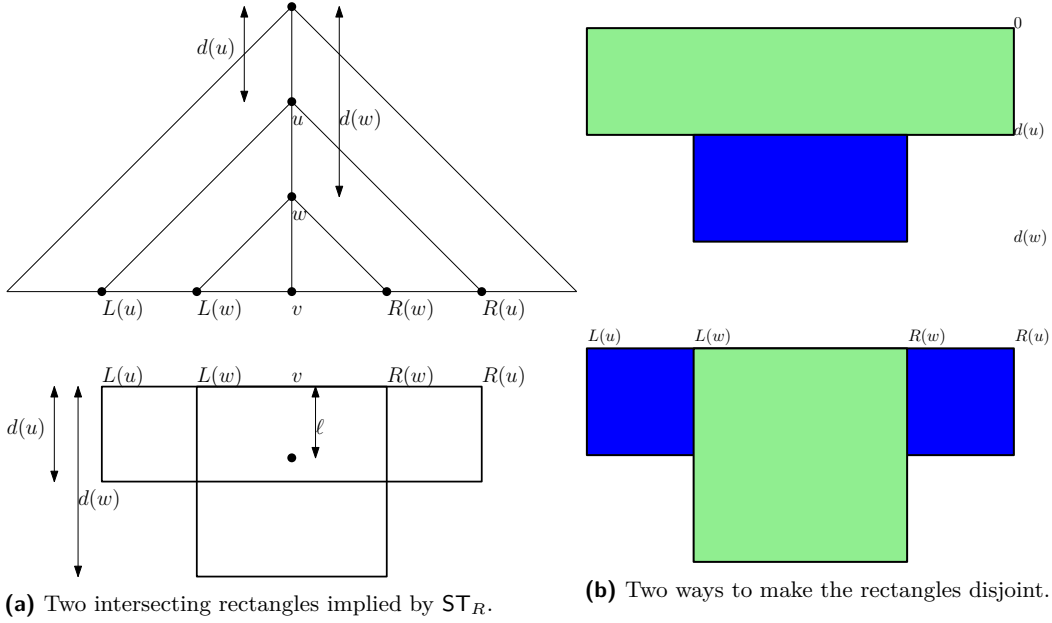
468 Now consider constructing the 2D rectangle stabbing data structure for reporting  
 469 (resp. counting) for these  $n$  rectangles, and then ask the query for a point  $(L(v), \ell)$ , where  $v$   
 470 is the leaf node reached from the root by conceptually reading  $S_j\$_j$ . The data structure will  
 471 report (resp. count) all of the suffixes of  $S_i$ , for  $i \in [1, k]$ , of length at least  $\ell$  that are also  
 472 prefixes of  $S_j$ . Unfortunately, such a solution differs from the expected results of  $\text{Report}(i, \ell)$   
 473 and  $\text{Count}(i, \ell)$  in the following two ways:

- 474 1. Instead of finding all  $j \in [1, k]$  such that  $\text{SPL}_{i,j} \geq \ell$  for a given  $i$ , we find all such  $i \in [1, k]$   
 475 for a given  $j$ . This issue is addressed by Observation 1, which states that  $\text{Report}(i, \ell)$   
 476 and  $\text{Count}(i, \ell)$  reduce trivially to the problems considered here, denoted by  $\text{Report}^r(i, \ell)$   
 477 and  $\text{Count}^r(i, \ell)$ , respectively (recall that the  $r$  superscript refers to reversing the input  
 478 strings);
- 479 2. If there are multiple prefixes of  $S_j$  of length at least  $\ell$  that are also suffixes of  $S_i$ , then  
 480 we will report (resp. count) each of them leading to double reporting (resp. counting).  
 481 Although one may actually be interested in reporting or counting those multiple suffix-  
 482 prefixes, in this paper, we are only interested in the *longest* ones. We address this issue  
 483 by modifying the rectangles before the construction.

484 As mentioned earlier the first issue is resolved by Observation 1. To solve the second  
 485 issue, we have to make the set of rectangles, for a single  $i \in [1, k]$ , pairwise disjoint while  
 486 leaving their union unchanged. Notice that two such non-disjoint rectangles must come from  
 487 a pair of nodes  $u$  and  $w$  in an ancestor-descendant relationship. An easy solution is to take,  
 488 for every node  $w$  which has an outgoing edge with label  $\$_i$ , its closest ancestor  $u$  which  
 489 also has an outgoing edge with label  $\$_i$ , and change the  $[L(w), R(w)] \times [0, d(w)]$  rectangle  
 490 into  $[L(w), R(w)] \times [d(u) + 1, d(w)]$ ; inspect Figure 6. Since the part  $[L(w), R(w)] \times [0, d(u)]$   
 491 is already contained in  $[L(u), R(u)] \times [0, d(u)]$  the union remains unchanged, and since  
 492  $u$  is the closest such ancestor, the other rectangles (for this  $i$ ) cannot have a nonempty  
 493 intersection with the newly obtained one (the intersection with the ones coming from the  
 494 descendants of  $w$  is empty after the modification of those rectangles). We can perform these

---

<sup>2</sup>  $[L(u), R(u)]$  is also known as the suffix array interval of node  $u$ .



■ **Figure 6** On the bottom left part, the rectangles obtained from two nodes  $u$  and  $w$  of  $ST_R$  (top left), both having an outgoing edge with label  $\$i$ , forming a suffix-prefix match of  $S_i$  and  $S_j$  for node  $v$  reached by reading  $S_j\$j$  from the root. The rectangles have a nonempty intersection. To avoid double reporting (or double counting), we make the rectangles disjoint while leaving their union unchanged. We can do this (by taking the intersection *once*) in two ways (on the right): a simple one (top) or a more complicated one (bottom), which allows us to efficiently output  $SPL_{i,j}$ .

495 modifications with a single DFS traversal with  $k$  stacks of nodes on the path from the root  
 496 to the currently processed node, which has an outgoing edge with label  $\$i$ ,  $i \in [1, k]$ . A more  
 497 complicated solution is obtained by replacing the two rectangles  $[L(u), R(u)] \times [0, d(u)]$  and  
 498  $[L(w), R(w)] \times [0, d(w)]$  with three rectangles:  $[L(u), L(w) - 1] \times [0, d(u)]$ ,  $[L(w), R(w)] \times$   
 499  $[0, d(w)]$  and  $[R(w) + 1, R(u)] \times [0, d(u)]$ ; inspect Figure 6. Unlike the previous construction, a  
 500 single rectangle can be spliced into smaller ones many times (a node can be a direct ancestor  
 501 of many other nodes); at the same time a single rectangle can splice only its direct ancestor,  
 502 hence the number of rectangles obtained this way is bounded from above by  $2n$ . This set of  
 503 modified intervals can be obtained similarly: in a DFS traversal, when a node which has  
 504 an outgoing edge with label  $\$i$  is reached, we access its closest ancestor, which also has an  
 505 outgoing edge with label  $\$i$ , and splice its rectangle. As such descendants of a node are  
 506 visited from left to right, we always know which part of the rectangle will be spliced next,  
 507 hence each such splice takes  $\mathcal{O}(1)$  time leading to computing  $\mathcal{O}(n)$  such modified rectangles  
 508 in  $\mathcal{O}(n)$  total time.

509 In order to finalize the construction of our data structure, we compute the set of modified  
 510 rectangles of one of the two types described above, and construct for them the 2D rectangle  
 511 stabbing data structures for reporting (Lemma 16) and counting (Lemma 17).

512 **Querying** To answer a  $\text{Report}^r(j, \ell)$  or a  $\text{Count}^r(j, \ell)$  query, we simply ask the corresponding  
 513 2D rectangle stabbing data structure for the point  $(L(v), \ell) = (R(v), \ell)$ , where  $v$  is the node  
 514 reached in  $ST_R$  from the root by conceptually reading  $S_j\$j$ . In case of a reporting query,  
 515 the data structure returns a set of rectangles  $[x, y] \times [\ell_1, \ell_2]$  labeled with distinct values  
 516  $i \in [1, k]$ . We can simply report the set of these  $i$  values. In case of a counting query, the

517 result is simply an integer which we output. The two constructions of modified rectangles  
 518 have additional nice properties however – each value  $i$  is associated with a value  $\ell_2$ . In case  
 519 of the first construction, this  $\ell_2$  is the length of the shortest suffix of  $S_i$  which is also a prefix  
 520 of  $S_j$  of length at least  $\ell$ ; in case of the second construction,  $\ell_2$  is the length of the longest  
 521 such suffix, that is  $\ell_2 = \text{SPL}_{i,j}$ .

522 **Correctness** The correctness of the algorithm follows by the fact that point  $(L(v), \ell) =$   
 523  $(R(v), \ell)$  is enclosed by a rectangle  $[L(u), R(u)] \times [0, d(u)]$  if and only if  $S_j\$j$  has a prefix  
 524 of length at least  $\ell$  that is also a suffix of  $S_i$ ; and by the fact that the set of rectangles  
 525 originating from a single  $i$  are made pairwise disjoint while their union remains unchanged.

526 We have thus arrived at the following lemma.

527 **► Lemma 18.** *For any dictionary of  $k$  strings of total length  $n$  over an integer alphabet of*  
 528 *size  $\sigma \leq n^{\mathcal{O}(1)}$ , we can construct a data structure of  $\mathcal{O}(n)$  words of space answering: (i)*  
 529  *$\text{Report}^r(j, \ell)$  queries in  $\mathcal{O}(\log n / \log \log n + f)$  time, where  $f$  is the size of the output; and*  
 530 *(ii)  $\text{Count}^r(j, \ell)$  queries in  $\mathcal{O}(\log n / \log \log n)$  time.*

531 By combining Lemma 18 with Observation 1 we obtain the main result of this section.

532 **► Theorem 19.** *For any dictionary of  $k$  strings of total length  $n$  over an integer alphabet*  
 533 *of size  $\sigma \leq n^{\mathcal{O}(1)}$ , we can construct a data structure of  $\mathcal{O}(n)$  words of space answering:*  
 534 *(i)  $\text{Report}(i, \ell)$  queries in  $\mathcal{O}(\log n / \log \log n + f)$  time, where  $f$  is the output size; and (ii)*  
 535  *$\text{Count}(i, \ell)$  queries in  $\mathcal{O}(\log n / \log \log n)$  time.*

536 Let us remark that the construction time for our data structures, excluding the imple-  
 537 mentation of the data structures underlying Lemmas 16 and 17, is  $\mathcal{O}(n)$ . Unfortunately,  
 538 the construction time of the latter data structures (Lemmas 16 and 17) is not mentioned  
 539 in [28, 45]. However, by using the construction from [15], we obtain  $\mathcal{O}(n \log n)$  construction  
 540 time,  $\mathcal{O}(n)$  words of space,  $\mathcal{O}(\log n + f)$  time for reporting, and  $\mathcal{O}(\log n)$  time for counting.

541 **► Theorem 20.** *For any dictionary of  $k$  strings of total length  $n$  over an integer alphabet*  
 542 *of size  $\sigma \leq n^{\mathcal{O}(1)}$ , we can construct a data structure of  $\mathcal{O}(n)$  words of space answering: (i)*  
 543  *$\text{Report}(i, \ell)$  queries in  $\mathcal{O}(\log n + f)$  time, where  $f$  is the output size; and (ii)  $\text{Count}(i, \ell)$*   
 544 *queries in  $\mathcal{O}(\log n)$  time. The data structure construction time is  $\mathcal{O}(n \log n)$ .*

545 Let us also remark that  $\text{Report}(i, 0)$  (with the second construction of disjoint rectangles)  
 546 actually answers any  $\text{One-to-All}(i)$  query within the same asymptotic time:  $\mathcal{O}(\log n + f) =$   
 547  $\mathcal{O}(\log n + k) = \mathcal{O}(k)$  as  $k \geq \sqrt{n}$ . While the data structure for answering  $\text{Report}$  queries  
 548 occupies  $\mathcal{O}(n)$  words of space, like the data structure for  $\text{One-to-All}$  queries, the construction  
 549 time for the former is more expensive – and it is likely much slower in practice.

## 550 **6 Answering Top Queries**

551 Recall that a  $\text{Top}(i, K)$  query returns exactly  $K$  elements  $j$  for which  $\text{SPL}_{i,j}$  is the largest,  
 552 breaking ties arbitrarily. In case we are given an additional bound  $K' \leq k$  such that  $K \leq K'$   
 553 (e.g., we are only interested in finding  $\mathcal{O}(1)$  many such top elements), the obvious data  
 554 structure would be to store, for each  $i \in [1, k]$ , the sorted list of size  $K'$  of the best answers.  
 555 Such a data structure allows answering  $\text{Top}(i, K)$  queries, for  $K \leq K'$ , in the optimal  $\mathcal{O}(K)$   
 556 time, but it requires  $\mathcal{O}(kK')$  space, which for small  $K'$  may be  $\mathcal{O}(n)$ , but in general (i.e.,  
 557 when  $K' = k$ ) leads back to the  $\mathcal{O}(n + k^2)$ -time APSP algorithm. We show how to use our  
 558 results from Section 5 to answer  $\text{Top}(i, K)$  queries using  $\mathcal{O}(n)$  space without this  $K'$  bound.

559 Clearly, we can assume that  $K < k$ . We start by making the following crucial observation.

560 ► **Observation 21.** *For any  $\text{Top}(i, K)$  query, with  $K < k$ , there exists an integer  $\ell \in [0, n - 1]$   
561 such that  $\text{Count}(i, \ell + 1) \leq K < \text{Count}(i, \ell)$ .*

562 Using the results from Section 5, we can find such an  $\ell$  in  $\mathcal{O}(\log^2 n / \log \log n)$  time using  
563 binary search on  $\ell \in [0, n - 1]$  and the data structure for  $\text{Count}$  queries. Next we can simply  
564 compute  $\text{Report}(i, \ell + 1)$  to be left with only choosing the remaining  $(K - \text{Count}(i, \ell + 1))$   
565 elements out of all  $j \in [1, k]$  such that  $\text{SPL}_{i,j} = \ell$ . Unfortunately, there can be many such  
566 elements (even  $k$ ), and we do not want this to influence the query time. We have to report  
567 the remaining elements out of the ones such that  $\text{SPL}_{i,j} = \ell$  without computing or explicitly  
568 accessing all of them. Recall that, in  $\text{ST}_R$ , a list of elements  $i$  such that  $S_i$  has a suffix of  
569 length exactly  $\ell$  which is also a prefix of  $S_j$  can be accessed in  $\mathcal{O}(1)$  time after  $\mathcal{O}(n)$ -time  
570 preprocessing by finding the ancestor of the node reached by conceptually reading  $S_j \$_j$  at  
571 string depth  $\ell$  (using a WA query) and reading the first letters of its outgoing edges from left  
572 to right; since  $\$_1 < \dots < \$_k$  are smaller than any element of  $\Sigma$  those values form a sorted  
573 list. Analogously, to access the list of elements  $j$  such that  $S_i$  has a suffix of length exactly  $\ell$   
574 which is also a prefix of  $S_j$ , we simply use the symmetric data structure by Observation 1.

575 Unfortunately, this list may contain elements  $j$  such that  $\text{SPL}_{i,j} > \ell$ , and we do not  
576 want to report them again. This, however, can be fixed by maintaining a bitvector of size  $k$   
577 as an integral part of our data structure; for each element  $j \in \text{Report}(i, \ell + 1)$ , we set the  
578  $j$ th element of the bitvector to 1 in  $\mathcal{O}(\text{Count}(i, \ell + 1)) = \mathcal{O}(K)$  time. When accessing the  
579 elements of the sorted list one-by-one, we simply check if the element was already outputted  
580 using the bitvector in  $\mathcal{O}(1)$  time. In total, we can check up to  $K$  such elements, hence the  
581 total time of merging those two parts of the output is  $\mathcal{O}(K)$  (including the bitvector reset).  
582 We summarize the solution in Theorem 22, which is the main result of this section.

583 ► **Theorem 22.** *For any dictionary of  $k$  strings of total length  $n$  over an integer alphabet of  
584 size  $\sigma \leq n^{\mathcal{O}(1)}$ , we can construct a data structure of  $\mathcal{O}(n)$  words of space answering  $\text{Top}(i, K)$   
585 queries in  $\mathcal{O}(\log^2 n / \log \log n + K)$  time.*

586 **Proof.** We start the construction of the data structure by constructing the data structures  
587 for  $\text{Report}(i, \ell)$  and  $\text{Count}(i, \ell)$  using Theorem 19. We also construct a data structure to  
588 find the list of elements  $j$  such that  $S_i$  has a suffix-prefix match of length  $\ell$  with  $S_j$  in  $\mathcal{O}(1)$   
589 time using Lemmas 2 and 7 and Observation 1. Finally, we also maintain a bitvector of size  
590  $k = \mathcal{O}(n)$ . The space required by our data structure is  $\mathcal{O}(n)$  words.

591 Consider a  $\text{Top}(i, K)$  query. We ask  $\mathcal{O}(\log n)$   $\text{Count}$  queries and a single  $\text{Report}$  query in  
592  $\mathcal{O}(\log^2 n / \log \log n + K)$  total time, as the output is bounded by  $K$ . We index the  $\text{Report}$   
593 result in the bitvector. We find the list (without reading its content) of elements  $j$  such that  
594  $S_i$  has a suffix of length exactly  $\ell$  which is also a prefix of  $S_j$  in  $\mathcal{O}(1)$  time. Finally, we access  
595 and check at most  $K$  elements from the list in  $\mathcal{O}(K)$  total time.

596 The correctness of the algorithm follows by Observation 21 and Theorem 19. ◀

597 Similar to Section 5, the construction time for our data structure, excluding the imple-  
598 mentation of Theorem 19, is  $\mathcal{O}(n)$ . If instead of Theorem 19, we employ Theorem 20, we  
599 obtain  $\mathcal{O}(n \log n)$  construction time,  $\mathcal{O}(n)$  words of space, and  $\mathcal{O}(\log^2 n + K)$  query time.

600 ► **Theorem 23.** *For any dictionary of  $k$  strings of total length  $n$  over an integer alphabet of  
601 size  $\sigma \leq n^{\mathcal{O}(1)}$ , we can construct a data structure of  $\mathcal{O}(n)$  words of space answering  $\text{Top}(i, K)$   
602 queries in  $\mathcal{O}(\log^2 n + K)$  time. The data structure construction time is  $\mathcal{O}(n \log n)$ .*



603 ——— **References** ———

- 604 1 Alfred V. Aho and Margaret J. Corasick. Efficient string matching: An aid to bibliographic  
605 search. *Commun. ACM*, 18(6):333–340, 1975.
- 606 2 Stephen Alstrup, Jacob Holm, Kristian de Lichtenberg, and Mikkel Thorup. Minimizing  
607 diameters of dynamic trees. In Pierpaolo Degano, Roberto Gorrieri, and Alberto Marchetti-  
608 Spaccamela, editors, *Automata, Languages and Programming, 24th International Colloquium,*  
609 *ICALP'97, Bologna, Italy, 7-11 July 1997, Proceedings*, volume 1256 of *Lecture Notes in*  
610 *Computer Science*, pages 270–280. Springer, 1997.
- 611 3 Amihoud Amir, Panagiotis Charalampopoulos, Solon P. Pissis, and Jakub Radoszewski.  
612 Dynamic and internal longest common substring. *Algorithmica*, 82(12):3707–3743, 2020.
- 613 4 Golnaz Badkobeh, Panagiotis Charalampopoulos, Dmitry Kosolobov, and Solon P. Pissis.  
614 Internal shortest absent word queries in constant time and linear space. *Theor. Comput. Sci.*,  
615 922:271–282, 2022.
- 616 5 Carl Barton, Costas S. Iliopoulos, Solon P. Pissis, and William F. Smyth. Fast and simple  
617 computations using prefix tables under hamming and edit distance. In Jan Kratochvíl, Mirka  
618 Miller, and Dalibor Fronček, editors, *Combinatorial Algorithms - 25th International Workshop,*  
619 *IWOCA 2014, Duluth, MN, USA, October 15-17, 2014, Revised Selected Papers*, volume 8986  
620 of *Lecture Notes in Computer Science*, pages 49–61. Springer, 2014.
- 621 6 Djamal Belazzougui, Dmitry Kosolobov, Simon J. Puglisi, and Rajeev Raman. Weighted  
622 ancestors in suffix trees revisited. In Paweł Gawrychowski and Tatiana Starikovskaya, editors,  
623 *32nd Annual Symposium on Combinatorial Pattern Matching, CPM 2021, July 5-7, 2021,*  
624 *Wrocław, Poland*, volume 191 of *LIPIcs*, pages 8:1–8:15. Schloss Dagstuhl - Leibniz-Zentrum  
625 für Informatik, 2021.
- 626 7 Djamal Belazzougui and Gonzalo Navarro. Optimal lower and upper bounds for representing  
627 sequences. *ACM Trans. Algorithms*, 11(4):31:1–31:21, 2015.
- 628 8 Ilan Ben-Bassat and Benny Chor. String graph construction using incremental hashing.  
629 *Bioinform.*, 30(24):3515–3523, 2014.
- 630 9 Michael A. Bender and Martin Farach-Colton. The LCA problem revisited. In Gaston H.  
631 Gonnet, Daniel Panario, and Alfredo Viola, editors, *LATIN 2000: Theoretical Informatics, 4th*  
632 *Latin American Symposium, Punta del Este, Uruguay, April 10-14, 2000, Proceedings*, volume  
633 1776 of *Lecture Notes in Computer Science*, pages 88–94. Springer, 2000.
- 634 10 Philip Bille and Inge Li Gørtz. The tree inclusion problem: In linear space and faster. *ACM*  
635 *Trans. Algorithms*, 7(3):38:1–38:47, 2011.
- 636 11 Paola Bonizzoni, Gianluca Della Vedova, Yuri Pirola, Marco Previtali, and Raffaella Rizzi.  
637 FSG: fast string graph construction for de novo assembly. *J. Comput. Biol.*, 24(10):953–968,  
638 2017.
- 639 12 Bastien Cazaux and Eric Rivals. Hierarchical overlap graph. *Inf. Process. Lett.*, 155, 2020.
- 640 13 Panagiotis Charalampopoulos, Paweł Gawrychowski, Shay Mozes, and Oren Weimann. Almost  
641 optimal distance oracles for planar graphs. In Moses Charikar and Edith Cohen, editors,  
642 *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing, STOC*  
643 *2019, Phoenix, AZ, USA, June 23-26, 2019*, pages 138–151. ACM, 2019.
- 644 14 Panagiotis Charalampopoulos, Tomasz Kociumaka, Manal Mohamed, Jakub Radoszewski,  
645 Wojciech Rytter, and Tomasz Walen. Internal dictionary matching. *Algorithmica*, 83(7):2142–  
646 2169, 2021.
- 647 15 Bernard Chazelle. A functional approach to data structures and its use in multidimensional  
648 searching. *SIAM J. Comput.*, 17(3):427–462, 1988.
- 649 16 Shiri Chechik. Approximate distance oracles with constant query time. In David B. Shmoys,  
650 editor, *Symposium on Theory of Computing, STOC 2014, New York, NY, USA, May 31 -*  
651 *June 03, 2014*, pages 654–663. ACM, 2014.
- 652 17 Shiri Chechik. Approximate distance oracles with improved bounds. In Rocco A. Servedio  
653 and Ronitt Rubinfeld, editors, *Proceedings of the Forty-Seventh Annual ACM on Symposium*

- 654 on *Theory of Computing, STOC 2015, Portland, OR, USA, June 14-17, 2015*, pages 1–10.  
 655 ACM, 2015.
- 656 **18** Nicola Cotumaccio, Giovanna D’Agostino, Alberto Policriti, and Nicola Prezza. Co-  
 657 lexicographically ordering automata and regular languages. part I. *CoRR*, abs/2208.04931,  
 658 2022.
- 659 **19** Maxime Crochemore, Christophe Hancart, and Thierry Lecroq. *Algorithms on strings*. Cam-  
 660 bridge University Press, 2007.
- 661 **20** Shiri Dori and Gad M. Landau. Construction of aho corasick automaton in linear time for  
 662 integer alphabets. *Inf. Process. Lett.*, 98(2):66–72, 2006.
- 663 **21** Herbert Edelsbrunner and Mark H. Overmars. On the equivalence of some rectangle problems.  
 664 *Inf. Process. Lett.*, 14(3):124–127, 1982.
- 665 **22** Martin Farach. Optimal suffix tree construction with large alphabets. In *38th Annual*  
 666 *Symposium on Foundations of Computer Science, FOCS ’97, Miami Beach, Florida, USA,*  
 667 *October 19-22, 1997*, pages 137–143. IEEE Computer Society, 1997.
- 668 **23** Martin Farach and S. Muthukrishnan. Perfect hashing for strings: Formalization and algorithms.  
 669 In Daniel S. Hirschberg and Eugene W. Myers, editors, *Combinatorial Pattern Matching, 7th*  
 670 *Annual Symposium, CPM 96, Laguna Beach, California, USA, June 10-12, 1996, Proceedings*,  
 671 volume 1075 of *Lecture Notes in Computer Science*, pages 130–140. Springer, 1996.
- 672 **24** Pawel Gawrychowski, Moshe Lewenstein, and Patrick K. Nicholson. Weighted ancestors in  
 673 suffix trees. In Andreas S. Schulz and Dorothea Wagner, editors, *Algorithms - ESA 2014*  
 674 *- 22th Annual European Symposium, Wroclaw, Poland, September 8-10, 2014. Proceedings*,  
 675 volume 8737 of *Lecture Notes in Computer Science*, pages 455–466. Springer, 2014.
- 676 **25** Giorgio Gonnella and Stefan Kurtz. Readjoinder: a fast and memory efficient string graph-based  
 677 sequence assembler. *BMC Bioinform.*, 13:82, 2012.
- 678 **26** Dan Gusfield. *Algorithms on Strings, Trees, and Sequences - Computer Science and Computa-*  
 679 *tional Biology*. Cambridge University Press, 1997.
- 680 **27** Dan Gusfield, Gad M. Landau, and Baruch Schieber. An efficient algorithm for the all pairs  
 681 suffix-prefix problem. *Inf. Process. Lett.*, 41(4):181–185, 1992.
- 682 **28** Joseph F. JáJá, Christian Worm Mortensen, and Qingmin Shi. Space-efficient and fast  
 683 algorithms for multidimensional dominance reporting and counting. In Rudolf Fleischer  
 684 and Gerhard Trippen, editors, *Algorithms and Computation, 15th International Symposium,*  
 685 *ISAAC 2004, Hong Kong, China, December 20-22, 2004, Proceedings*, volume 3341 of *Lecture*  
 686 *Notes in Computer Science*, pages 558–568. Springer, 2004.
- 687 **29** Shahbaz Khan. Optimal construction of hierarchical overlap graphs. In Pawel Gawrychowski  
 688 and Tatiana Starikovskaya, editors, *32nd Annual Symposium on Combinatorial Pattern*  
 689 *Matching, CPM 2021, July 5-7, 2021, Wroclaw, Poland*, volume 191 of *LIPICs*, pages 17:1–  
 690 17:11. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2021.
- 691 **30** Tomasz Kociumaka. Efficient data structures for internal queries in texts. *PhD thesis,*  
 692 *University of Warsaw, October 2018.*, 2018.
- 693 **31** Tomasz Kociumaka, Jakub Radoszewski, Wojciech Rytter, and Tomasz Walen. Internal  
 694 pattern matching queries in a text and applications. In Piotr Indyk, editor, *Proceedings of*  
 695 *the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2015, San*  
 696 *Diego, CA, USA, January 4-6, 2015*, pages 532–551. SIAM, 2015.
- 697 **32** Gregory Kucherov and Dekel Tsur. Improved filters for the approximate suffix-prefix overlap  
 698 problem. In Edleno Silva de Moura and Maxime Crochemore, editors, *String Processing*  
 699 *and Information Retrieval - 21st International Symposium, SPIRE 2014, Ouro Preto, Brazil,*  
 700 *October 20-22, 2014. Proceedings*, volume 8799 of *Lecture Notes in Computer Science*, pages  
 701 139–148. Springer, 2014.
- 702 **33** Jihyuk Lim and Kunsoo Park. A fast algorithm for the all-pairs suffix-prefix problem. *Theor.*  
 703 *Comput. Sci.*, 698:14–24, 2017.
- 704 **34** Grigorios Loukides and Solon P. Pissis. All-pairs suffix/prefix in optimal time using Aho-  
 705 Corasick space. *Inf. Process. Lett.*, 178:106275, 2022.

- 706 **35** Felipe A. Louza, Simon Gog, Leandro Zanotto, Guido Araujo, and Guilherme P. Telles. Parallel  
707 computation for the all-pairs suffix-prefix problem. In Shunsuke Inenaga, Kunihiko Sadakane,  
708 and Tetsuya Sakai, editors, *String Processing and Information Retrieval - 23rd International*  
709 *Symposium, SPIRE 2016, Beppu, Japan, October 18-20, 2016, Proceedings*, volume 9954 of  
710 *Lecture Notes in Computer Science*, pages 122–132, 2016.
- 711 **36** Udi Manber and Eugene W. Myers. Suffix arrays: A new method for on-line string searches.  
712 *SIAM J. Comput.*, 22(5):935–948, 1993.
- 713 **37** Eugene W. Myers. The fragment assembly string graph. *Bioinformatics*, 21(suppl\_2):ii79–ii85,  
714 09 2005.
- 715 **38** Gonzalo Navarro. *Compact Data Structures - A Practical Approach*. Cambridge University  
716 Press, 2016.
- 717 **39** Enno Ohlebusch and Simon Gog. Efficient algorithms for the all-pairs suffix-prefix problem  
718 and the all-pairs substring-prefix problem. *Inf. Process. Lett.*, 110(3):123–128, 2010.
- 719 **40** Sangsoo Park, Sung Gwan Park, Bastien Cazaux, Kunsoo Park, and Eric Rivals. A linear  
720 time algorithm for constructing hierarchical overlap graphs. In Pawel Gawrychowski and  
721 Tatiana Starikovskaya, editors, *32nd Annual Symposium on Combinatorial Pattern Matching,*  
722 *CPM 2021, July 5-7, 2021, Wroclaw, Poland*, volume 191 of *LIPICs*, pages 22:1–22:9. Schloss  
723 Dagstuhl - Leibniz-Zentrum für Informatik, 2021.
- 724 **41** Mihai Patrascu and Liam Roditty. Distance oracles beyond the thorup-zwick bound. *SIAM J.*  
725 *Comput.*, 43(1):300–311, 2014.
- 726 **42** Maan Haj Rachid and Qutaibah Malluhi. A practical and scalable tool to find overlaps between  
727 sequences. *BioMed Res. Int.*, 2015(905261), 2015.
- 728 **43** Rajeev Raman, Venkatesh Raman, and S. Srinivasa Rao. Succinct indexable dictionaries with  
729 applications to encoding k-ary trees and multisets. In David Eppstein, editor, *Proceedings of*  
730 *the Thirteenth Annual ACM-SIAM Symposium on Discrete Algorithms, January 6-8, 2002,*  
731 *San Francisco, CA, USA*, pages 233–242. ACM/SIAM, 2002.
- 732 **44** Kim R. Rasmussen, Jens Stoye, and Eugene W. Myers. Efficient  $q$ -gram filters for finding all  
733  $\epsilon$ -matches over a given length. *J. Comput. Biol.*, 13(2):296–308, 2006.
- 734 **45** Qingmin Shi and Joseph F. JáJá. Novel transformation techniques using  $q$ -heaps with  
735 applications to computational geometry. *SIAM J. Comput.*, 34(6):1474–1492, 2005.
- 736 **46** Jared T. Simpson and Richard Durbin. Efficient construction of an assembly string graph  
737 using the fm-index. *Bioinform.*, 26(12):367–373, 2010.
- 738 **47** Sharma V. Thankachan, Chaitanya Aluru, Sriram P. Chockalingam, and Srinivas Aluru.  
739 Algorithmic framework for approximate matching under bounded edits with applications to  
740 sequence analysis. In Benjamin J. Raphael, editor, *Research in Computational Molecular*  
741 *Biology - 22nd Annual International Conference, RECOMB 2018, Paris, France, April 21-*  
742 *24, 2018, Proceedings*, volume 10812 of *Lecture Notes in Computer Science*, pages 211–224.  
743 Springer, 2018.
- 744 **48** Mikkel Thorup and Uri Zwick. Approximate distance oracles. *J. ACM*, 52(1):1–24, 2005.
- 745 **49** William H. A. Tustumi, Simon Gog, Guilherme P. Telles, and Felipe A. Louza. An improved  
746 algorithm for the all-pairs suffix-prefix problem. *J. Discrete Algorithms*, 37:34–43, 2016.
- 747 **50** Esko Ukkonen. A linear-time algorithm for finding approximate shortest common superstrings.  
748 *Algorithmica*, 5(3):313–323, 1990.
- 749 **51** Esko Ukkonen. On-line construction of suffix trees. *Algorithmica*, 14(3):249–260, 1995.
- 750 **52** Niko Välimäki, Susana Ladra, and Veli Mäkinen. Approximate all-pairs suffix/prefix overlaps.  
751 *Inf. Comput.*, 213:49–58, 2012.
- 752 **53** Peter Weiner. Linear pattern matching algorithms. In *14th Annual Symposium on Switching and*  
753 *Automata Theory, Iowa City, Iowa, USA, October 15-17, 1973*, pages 1–11. IEEE Computer  
754 Society, 1973.