

1 Practical Performance of Space Efficient Data 2 Structures for Longest Common Extensions

3 **Patrick Dinklage**

4 Department of Computer Science, Technical University of Dortmund
5 Dortmund, Germany
6 patrick.dinklage@tu-dortmund.de

7 **Johannes Fischer**

8 Department of Computer Science, Technical University of Dortmund
9 Dortmund, Germany
10 johannes.fischer@cs.tu-dortmund.de

11 **Alexander Herlez**

12 Department of Computer Science, Technical University of Dortmund
13 Dortmund, Germany
14 alexander.herlez@tu-dortmund.de

15 **Tomasz Kociumaka**

16 Department of Computer Science, Bar Ilan University
17 Ramat Gan, Israel
18 kociumaka@mimuw.edu.pl

19 **Florian Kurpicz**

20 Department of Computer Science, Technical University of Dortmund
21 Dortmund, Germany
22 florian.kurpicz@tu-dortmund.de

23 — Abstract —

24 For a text $T[1, n]$, a Longest Common Extension (LCE) query $\text{lce}_T(i, j)$ asks for the length of
25 the longest common prefix of the suffixes $T[i, n]$ and $T[j, n]$ identified by their starting positions
26 $1 \leq i, j \leq n$. A classic problem in stringology asks to preprocess a static text $T[1, n]$ over an alphabet
27 of size σ so that LCE queries can be efficiently answered on-line. Since its introduction in 1980's, this
28 problem has found numerous applications: in suffix sorting, edit distance computation, approximate
29 pattern matching, regularities finding, string mining, and many more. Text-book solutions offer
30 $O(n)$ preprocessing time and $O(1)$ query time, but they employ memory-heavy data structures, such
31 as suffix arrays, in practice several times bigger than the text itself.

32 Very recently, more space efficient solutions using $O(n \log \sigma)$ bits of total space or even only
33 $O(\log n)$ bits of extra space have been proposed: string synchronizing sets [Kempa and Kociumaka,
34 STOC'19, and Birenzwe et al., SODA'20] and in-place fingerprinting [Prezza, SODA'18]. The goal
35 of this article is to present well-engineered implementations of these new solutions and study their
36 practicality on a commonly agreed text corpus. We show that both perform extremely well in practice,
37 with space consumption of only around 10% of the input size for string synchronizing sets (around
38 20% for highly repetitive texts), and essentially no extra space for fingerprinting. Interestingly,
39 our experiments also show that both solutions become much faster than naive scanning even for
40 finding common prefixes of moderate length, contradicting a common belief that sophisticated data
41 structures for LCE queries are not competitive with naive approaches [Ilie and Tinta, SPIRE'09].

42 **2012 ACM Subject Classification** Theory of computation – Design and analysis of algorithms –
43 Data structures design and analysis – Pattern matching

44 **Keywords and phrases** text indexing, longest common prefix, space efficient data structures

45 **Funding** *Florian Kurpicz*: This work was supported by the German Research Foundation (DFG),
46 priority programme “Algorithms for Big Data” (SPP 1736).

47 **1 Introduction and Related Work**

48 The longest common extension (LCE) problem is formally defined as follows:

49 **Given:** A text $T[1, n]$ with n symbols from an alphabet of size σ .

50 **Construct:** A data structure that allows subsequent on-line queries of the form

$$51 \text{ lce}_T(i, j) := \max\{\ell \geq 0 : T[i, i + \ell] = T[j, j + \ell]\}.$$

52 This is a fundamental problem in stringology, with numerous applications in (sparse) suffix
53 sorting [12, 19, 24], edit distance computation [35, 37], approximate pattern matching [2, 22, 36],
54 regularities finding [4, 5, 13], string mining [21], and many more. The problem has a classic
55 solution with $O(1)$ -time queries based either on the suffix tree of T combined with lowest
56 common ancestor queries [15, 26, 37], or using the inverse suffix array of T and the longest
57 common prefix array with a data structure for range minimum queries [18, 30]. The latter
58 variant is more space efficient in practice but still several times larger than the text itself.

59 Recently, two notable methods appeared that also partly address the LCE problem:
60 Prezza’s *in-place fingerprinting* data structure [44] using a clever replacement of the original
61 text by well-chosen Karp–Rabin fingerprints [31], and data structures relying on local
62 consistency techniques [33, 39, 46] to construct so-called *string synchronizing sets* (SSS) [12, 32].

63 Although the LCE problem is also of high practical importance, we are only aware of
64 a few experimental papers on the LCE problem [18, 29]. However, these studies are now
65 somewhat dated, and naturally none of them includes the recent theoretical advances on the
66 problem. As a result, practitioners (e.g., in bioinformatics) still use those old (slow and/or
67 memory-heavy) data structures (see, e.g., [43, 45, 47, 50]), thereby limiting themselves to
68 problem sizes much smaller than they could actually handle.

69 **Our Contributions and Outline.** The goal of this paper is to engineer and evaluate LCE
70 data structures based on the recent theoretical advances in this field [12, 32, 44]: in-place
71 fingerprinting and string synchronizing sets. We first describe the theory behind these
72 two approaches in Sect. 3, after having introduced some general background in Sect. 2.
73 Sect. 4 describes further details of the implementations, including a simple but practical
74 data structure for the well-known predecessor/successor queries [49]. Finally, in Sect. 5, we
75 evaluate our implementations on a well-established benchmarking set, showing that (1) our
76 implementations based on SSS are always among the fastest to answer queries and (2) in-place
77 fingerprinting is very fast in practice and useful to answer queries that have long results.

78 **Further Related Work.** Compressed suffix trees with LCA functionality [20] can also be
79 used for LCE queries, but they are certainly too powerful (and hence too space consuming)
80 for the singular problem considered in this paper; this is confirmed experimentally in Sect. 5.

81 Time-space tradeoffs for LCE queries were considered by many authors [7, 9, 10, 12, 24, 48].
82 Kosolobov [34] showed that the product of extra space (in bits) and LCE query time must
83 be at least $\Omega(n \log n)$ under certain assumptions. Furthermore, there are data structures for
84 LCE queries in compressed [6, 8, 27, 28, 42, 48] and dynamic texts [1, 23, 39, 42].

85 **2 Preliminaries**

86 **2.1 Karp–Rabin Fingerprints**

87 The goal of Karp–Rabin fingerprints [31] is to hash substrings to small integers in order to
88 achieve fast comparisons for equality. To construct Karp–Rabin fingerprints, we choose a

89 random prime number $q = \Theta(n^c)$ for some constant $c > 1$. The Karp–Rabin fingerprint
 90 $\phi(x, y)$ of the substring $T[x, y]$ is then defined as

$$91 \quad \phi(x, y) = \left(\sum_{z=x}^y T[z] \cdot \sigma^{y-z} \right) \bmod q. \quad (1)$$

92 Observe that the Karp–Rabin fingerprints of matching substrings are equal, i.e., for every
 93 integer $\ell \geq 0$, if $T[x, x + \ell] = T[y, y + \ell]$, then $\phi(x, x + \ell) = \phi(y, y + \ell)$. On the other hand,
 94 if two substrings do not match, their fingerprints will be different with high probability.
 95 Specifically, if $T[x, x + \ell] \neq T[y, y + \ell]$, then $\text{Prob}[\phi(x, x + \ell) = \phi(y, y + \ell)] = O\left(\frac{\ell \log \sigma}{n^c}\right)$. Thus,
 96 by choosing large enough constant $c > 1$, we can control the probability of false positives
 97 when comparing fingerprints instead of the underlying substrings.

98 2.2 Static Successor Data Structures

99 Let U be a universe of u subsequent integers and let $S[0, n - 1]$ be a sequence of n integers
 100 from U sorted in ascending order. For any $x \in U$, we call $\text{succ}_S(x) := \min\{y \in S \mid y \geq x\}$
 101 the *successor* of x in S , i.e., the smallest value in S larger than or equal to x . Without loss
 102 of generality, we assume that $x > S[0]$ (otherwise, its successor is $S[0]$) and $x \leq S[n - 1]$
 103 (otherwise, it has no successor). Furthermore, in the following, we are interested only in the
 104 *position* i of the successor such that $S[i] = \text{succ}_S(x)$.

105 The most straightforward ways to find $\text{succ}_S(x)$ is by doing a linear search in time $O(n)$
 106 or, since S is sorted, by doing a binary search in time $O(\log n)$. Both require $O(\log n)$ bits
 107 of extra space for keeping track of the current position or the search interval and pivot,
 108 respectively. Alternatively, we can construct a bit vector B_S of size $u' = S[n - 1] - S[0] + 1$
 109 bits, in which we set $B_S[x - S[0]] := 1$ if and only if $x \in S$. Enhancing B_S to support
 110 constant-time rank queries (see [41, p. 75]) then allows for constant-time successor queries,
 111 because $\text{succ}_S(x) = \text{rank}_1(B_S, x - S[0])$. However, this method requires $u' + o(u')$ bits of
 112 memory and only works for static sets. Matsuoka et al. [38, Corollary 1] showed how to
 113 achieve the same space bounds in a dynamic setting, raising the query time to $O(\log \log u)$.

114 One can also interpret the numbers from S as binary strings and store them in a trie;
 115 some nodes of that trie can then be sampled to speed up the successor search. Such
 116 *universe-based sampling* has already been used in the *van Emde Boas tree* [49], which answers
 117 successor queries in time $O(\log \log u)$ and requires $O(n \log u)$ bits of memory using hashing.
 118 Dementiev et al. [14] described a careful implementation of this idea for 32-bit universes
 119 called *STree*. However, the *STree* is designed for *dynamic* sequences and contains components
 120 to support insertion and deletion into the set, which are not needed in the static case.

121 3 LCE Data Structures

122 3.1 In-Place Fingerprinting

123 Prezza [44] showed how to store the Karp–Rabin fingerprints from Sect. 2.1 succinctly so
 124 that the fingerprints *replace* the original text characters, but the original characters can still
 125 be recovered.

126 To facilitate explanation, let us assume a byte-alphabet ($\sigma = 256$) and a computer with
 127 word size $w := 64$ bits. We group characters from T into blocks of size $\tau := \Theta(w / \log \sigma)$;
 128 in our case, $\tau = 8$. Call the resulting array $B[1, n/\tau]$, with $B[i] := T[(i - 1)\tau + 1, i\tau]$
 129 for all $i = 1, \dots, n/\tau$ (assume for simplicity that τ divides n). Note that τ characters fit

130 into a computer word, and can hence be transferred to/from memory and arithmetically
 131 manipulated in constant time. We also choose a random prime q such that $\frac{1}{2}\sigma^\tau \leq q < \sigma^\tau$.
 132 We then compute the fingerprints $\phi(1, i\tau)$ for all $i = 1, \dots, n/\tau$. (The original description
 133 actually computes the values $(\phi(1, i\tau) + \bar{s}) \bmod q$, where $0 \leq \bar{s} < q$ is a random seed, which
 134 we omit in the following.) Since $D[i] := \lfloor B[i]/q \rfloor \in \{0, 1\}$, all we need to recover the full
 135 contents of block $B[i]$ are the two fingerprints $\phi(1, (i-1)\tau)$ and $\phi(1, i\tau)$, and the value $D[i]$.
 136 Prezza then shows how to choose q such that w.h.p. the fingerprints $\phi(1, i\tau)$ all have their
 137 most significant bit (MSB) set to 0; thus, the array D can be stored instead of those MSBs.
 138 (If this doesn't hold, one can recompute the data structure with a different value of q until
 139 it holds.) This also allows us to compute the fingerprints for *arbitrary* positions in T (not
 140 necessarily aligned with the block boundaries), also in $O(1)$ time. Prezza shows that with the
 141 above choices of q and τ , the fingerprints are collision free with high probability $1 - n^{-\Omega(1)}$.

142 The advantage of replacing the original text characters with fingerprints is that now,
 143 arbitrarily long substrings of the text can be tested for equality in constant time, using the
 144 fingerprints of those substrings. This can be applied for longest common extension queries
 145 as follows. To compute $\ell := \text{lce}_T(i, j)$ for any $1 \leq i, j \leq n$, we do an exponential search
 146 by comparing $\phi(i, i + 2^k)$ with $\phi(j, j + 2^k)$ for increasing exponents k until the fingerprints
 147 mismatch; the actual position of the first mismatch between $T[i, n]$ and $T[j, n]$ is then found
 148 by a further binary search on an interval of size $O(\ell)$. The whole process takes $O(\log \ell)$ time,
 149 assuming that we have access to a precomputed table of all necessary powers of σ modulo q ,
 150 which adds another $O(w \log n)$ bits of space (negligible in practice).

151 3.2 String Synchronizing Sets

152 The idea behind a *string synchronizing set* (SSS for short) is to designate a (hopefully small)
 153 set of positions from T such that this choice of positions is *locally consistent*, meaning that
 154 in sufficiently long matching substrings of T , the chosen positions are the same (relative to
 155 the beginnings of these substrings). We use the following simplified definition of SSS that is
 156 sufficient for our purposes:

157 ► **Definition 1.** For a positive integer $\tau \leq n/2$, the τ -synchronizing set of T is defined as

$$158 \quad S = \{i \in [1, n-2\tau+1] : \min\{\phi(j, j+\tau-1) : j \in [i, i+\tau]\} \in \{\phi(i, i+\tau-1), \phi(i+\tau, i+2\tau-1)\}\}.$$

159 The τ -synchronizing set satisfies the following *consistency property*: for all $i, j \in [1, n-2\tau+1]$,
 160 if $T[i, i+2\tau-1] = T[j, j+2\tau-1]$, then $i \in S \iff j \in S$. The original definition [32] also
 161 includes a *density property* that is necessary to guarantee a small SSS (of size $O(n/\tau)$) even
 162 for highly repetitive parts of T , which we do not consider in this paper, as it would make the
 163 data structure and query procedure significantly more complicated. (Omitting the density
 164 property can only make our data structure larger, but it remains correct).

165 We can use SSS for LCE data structures as follows. We first build a successor data
 166 structure on S . Let $n' := |S|$ denote the size of the SSS, and let $s_1 < s_2 < \dots < s_{n'}$
 167 denote the positions of S in increasing order. Define a new text T' of length n' by setting
 168 $T'[i] := T[s_i, s_i + 3\tau - 1]$ for all $1 \leq i \leq n'$. Then, build a data structure for constant-time
 169 LCE queries on T' , e.g., the RMQ-based solution mentioned in the introduction [18].

170 In order to answer an $\text{lce}_T(i, j)$ query for arbitrary $1 \leq i, j \leq n$, first compare at most
 171 $3\tau + 1$ characters from $T[i, i + 3\tau]$ and $T[j, j + 3\tau]$. If this comparison yields a mismatch, we
 172 are done. If not, compute $s_{i'} := \text{succ}_S(i)$, $s_{j'} := \text{succ}_S(j)$, and $\ell' := \text{lce}_{T'}(i', j')$. Then,

$$173 \quad \text{lce}_T(i, j) = s_{i'+\ell'} - i + \text{lce}_{T'}(s_{i'+\ell'}, s_{j'+\ell'}) , \quad (2)$$

174 and the latter lce_T value can again be computed naively as it is at most 3τ due to $T'[i'+\ell'] \neq$
 175 $T'[j'+\ell']$. Overall, we get $O(\tau)$ query time. (The original description [32] sets $\tau = \Theta(\log_\sigma n)$
 176 and uses the bit-vector approach from Sect. 2.2 for successor queries, as well as word-packing
 177 and bit-fiddling techniques, to achieve $O(1)$ query time, still in $O(n \log \sigma)$ bits of space.)

178 An advantage of this LCE data structure is that it naturally combines fast naive scanning
 179 for small LCE values (up to 3τ) with more sophisticated data structures guaranteeing a
 180 good worst-case performance. Taken together with the fact that the data structure is easily
 181 tunable (by adjusting τ), it is an excellent candidate for a practical LCE data structure.

182 4 Implementation

183 In this section, we describe our C++ implementations of the data structures from Sect. 3.
 184 The code is optimized for byte-alphabets ($\sigma = 256$) and can handle texts of arbitrary length,
 185 with the restriction for the algorithms described in Sect. 4.4 that the synchronizing set S
 186 can contain at most $|S| < 2^{32}$ elements.

187 4.1 A Simple but Fast Static Successor Data Structure

188 We first describe our implementation of a fast practical data structure for successor queries
 189 on a static sorted integer sequence $S[0, n-1]$ of length n over a universe U of size $u = 2^w$.
 190 To the best of our knowledge, there is no systematic study of such data structures. Our data
 191 structure is a simple index that combines universe-based sampling, binary search, and linear
 192 search to find the successor $\text{succ}_S(x)$ for any given $x \in U$.

193 Let $k < w$ be a trade-off parameter and let c be the number of elements from U that fit
 194 into a cache line. We build an index P over S such that we can, in constant time, look up
 195 an interval in S of length at most $2^k + 1$ in which $\text{succ}_S(x)$ is guaranteed to be contained. In
 196 this interval, we proceed with a binary search up to the point when the size of the current
 197 search interval becomes at most c . In this final small interval, which fully fits into a cache
 198 line, we look for $\text{succ}_S(x)$ using linear search. This successor query takes $O(k + c)$ time.

199 It remains to construct P . Let $\text{hi}_k(x) = \lfloor x/2^k \rfloor$. We define an array I_P of $2^{w-k} = u/2^k$
 200 intervals such that, for every non-negative integer $q < u/2^k$, the entry $I_P[q]$ is defined as

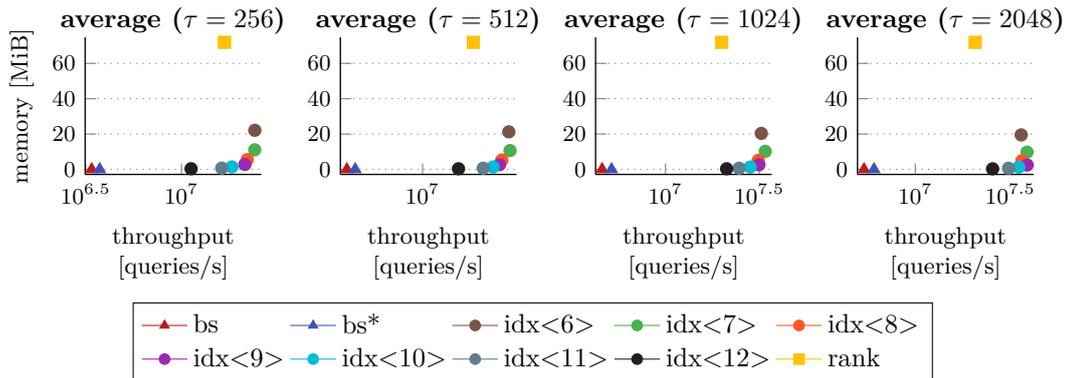
$$201 \quad \left[\max\{i \in [0, n-1] \mid \text{hi}_k(S[i]) < q\} + 1, \max\{j \in [0, n-2] \mid \text{hi}_k(S[j]) \leq q\} + 1 \right].$$

202
 203 If any maximum does not exist, we set the corresponding boundary to 0. As $\text{hi}_k(s) \leq u/2^k$
 204 for each $s \in S$, the right boundary of the last interval is always $n-1$. Informally, we split
 205 the universe U into 2^{w-k} intervals of size 2^k , find the corresponding boundaries in S , and
 206 store them in I_P with the right boundary shifted forward by one. The successor of x is then
 207 guaranteed to be contained in the interval $I_P[\text{hi}_k(x)]$ of S , which is of length at most $2^k + 1$.

208 Since, for $1 < q < u/2^k$, the right boundary of $I_P[q-1]$ equals the left boundary of
 209 $I_P[q]$, in our final index P , we only store the left boundaries and append $n-1$ (the right
 210 boundary of the last interval) at the end. Then, the interval of S in which the successor of
 211 x is contained can be expressed as $[P[\text{hi}_k(x)], P[\text{hi}_k(x) + 1]]$. We can construct P in time
 212 $O(n + u/2^k)$ with one scan over S and observe that P takes $(u/2^k + 1)\lceil \lg n \rceil$ bits of space.

213 Note that the space requirement is not optimal in theory and very wasteful if $u \gg n$.
 214 However, in our use case of τ -synchronizing sets, the universe is relatively small and gaps
 215 between two subsequent integers in S are of length at most τ , making our index practical.

216 In a preliminary experiment, we looked for the best configuration of the parameter k for
 217 our successor data structure. For this experiment, we constructed the string synchronizing



■ **Figure 1** Query throughput versus space usage for static successor data structures built on top of string synchronizing sets with various values of τ , averaged over all texts in the corpus (Tbl. 1).

218 sets S for our input texts (see Tbl. 1) and compared the following successor data structures:
 219 **1.** simple binary search (bs), **2.** binary search that switches to linear search (bs^*) for intervals
 220 of length at most c , **3.** the simple data structure from Sect. 2.2, using a bit vector supporting
 221 constant-time *rank* queries, and **4.** our index P for different parameters k ($idx\langle k \rangle$).

222 We measured the memory consumption of the successor data structures and the median
 223 query throughput for ten million successor queries over five iterations. The results in
 224 Figure 1 indicate that $k := 7$ yields the most beneficial trade-off between query throughput
 225 and memory consumption. Therefore, we used this configuration in further experiments.
 226 Appendix B shows a detailed breakdown of the results for all input texts.

227 4.2 Naive LCE Queries

228 Let us now turn our attention to LCE data structures. An **ultra naive** LCE query algorithm
 229 compares characters one by one until a mismatch is detected. This can be sped up significantly
 230 by using `uint_128` variables—nowadays supported by most CPUs—to compare blocks of
 231 16 characters at once. Once we have a mismatch, we compare the last block character by
 232 character. This algorithm is called **naive**.

233 4.3 In-Place Fingerprinting

234 We group 8 consecutive bytes of the input text into a block, reverse the order of the characters
 235 in each block by simply swapping the characters (to speed up arithmetic operations when
 236 querying the data structure), compute the fingerprint up to the end of the block using Eq. (1),
 237 and overwrite the block with this fingerprint. Reversing the order of characters has no
 238 significant affect the construction time, as this construction algorithm is still the fastest on
 239 all tested inputs (ignoring the naive approaches that do not require any construction at all),
 240 see Sect. 5. We use the smallest prime $q > 2^{63}$, which is $q = 2^{63} + 29$; this worked without
 241 any collision detected and with $MSB = 0$ for all stored fingerprints across all tested texts.
 242 We also use `uint_128` and the appropriate extended processor instructions for multiplying
 243 two 64-bit fingerprints.

244 Answering LCE queries using the procedure explained in Sect. 2.1 is done block-wise, i.e.,
 245 in steps by 8, 16, \dots characters. Within a block, instead of reconstructing and comparing
 246 fingerprints of short substrings, we first translate the fingerprint of the block back to its
 247 original contents (characters from T) and scan these characters naively. Because this is faster

248 for LCE values smaller than 256, we also do this at the start of every LCE query. That
 249 means that we only start the exponential search after scanning 8 blocks of 8 characters each.
 250 We do the same when the final binary search reaches an interval shorter than 256 characters.

251 4.4 String Synchronizing Sets

252 While constructing a synchronizing set S , we only need to keep the fingerprints of a window
 253 of size 2τ in memory. For reasons of speed, we do so by using a ring buffer of size $2^{\lceil \lg 2\tau \rceil}$:
 254 the power of two allows us to compute the positions in the buffer using fast bit operations.

255 When sorting the suffixes of T' , we first use sequential radix sort (`bingmann_msd_CI3_sb`)
 256 by Bingmann et al. [11] to reduce the alphabet size of T' (its characters are length- 3τ
 257 substrings of T). Then, we compute the suffix array using SAIS-LITE [40], because it is
 258 the fastest for integer alphabets [3]. According to Sect. 3.2, we now have to compute the
 259 LCP array of T' . However, since in the end we are interested only in the LCP values within
 260 the original text T , we deviate slightly from Eq. (2) and compute a *sparse* LCP array of T ,
 261 using the positions from S as indexed positions. For this, we rely on [32, Fact 5.1] that the
 262 lexicographic order of suffixes of T' coincides with the order of the corresponding suffixes
 263 of T . Apart from eliminating the need to map positions from T' back to T , this also saves
 264 us from performing the 3τ naive character comparisons at the end of the query in Eq. (2).

265 For the actual LCE computation for positions in S , we also build an RMQ data structure
 266 on the sparse LCP array; we apply a data structure by Ferrada and Navarro [16] for this
 267 purpose. We use a trick to speed up RMQ computations: since we know the length of the
 268 interval on which an RMQ is performed, we can make use of the fact that scanning small
 269 intervals is faster than an RMQ [29]. In our implementation, we only use the complicated
 270 RMQ machinery for intervals larger than 1,024; this value yielded the fastest query times.
 271 Smaller intervals are simply scanned for the minimum. We implemented two variants of LCE
 272 queries: **Prefer-short** corresponds to the description in Sect. 3.2: First, scan 3τ characters
 273 from T naively; if they are equal, compute the successors, and then apply Eq. (2) (with the
 274 modification described in the preceding paragraph). This variant should be used when *small*
 275 LCE values are expected, as the procedure is likely to stop already in the initial scan (before
 276 computing the successors). **Prefer-long**, on the other hand, is optimized for a setting where
 277 *large* LCE values are expected. Here, the important observation is that the initial naive
 278 scan could reach synchronized positions much earlier than after 3τ comparisons. From that
 279 point on, one could immediately resort to LCE computation on T' . We therefore swap the
 280 computation of the successors ($s_{i'}$ and $s_{j'}$) and the naive scan. The naive scan then only has
 281 to verify that $T[i, s_{i'}]$ and $T[j, s_{j'}]$ indeed match.

282 5 Experimental Evaluation

283 We tested the data structures from Sect. 4: **ultra_naive** and **naive** (Sect. 4.2), **our-**
 284 **rk** (fingerprinting from Sect. 4.3), and **sss _{τ}** (string synchronizing sets from Sect. 4.4
 285 with $\tau = 256, 512, 1024, 2048$) in the both variants (*pl* denotes the prefer-long vari-
 286 ant). Our implementations are available at <https://github.com/herlez/lce-test>. We
 287 also compared our implementations with existing implementations: **prezza-rk** (<https://github.com/nicolaprezza/rk-lce>) and data structures using the compressed suffix trees
 288 (CST) **sada** and **sct3** contained in SDSL [25] (<https://github.com/simongog/sdsl-lite>).
 289 Other existing implementations were excluded due to their much higher space consumption.
 290

291 Our experiments were conducted on a computer with two Intel Xeon E5-2640v4 CPUs
 292 (each with 10 cores at 2.4 GHz base frequency (3.4 GHz maximum turbo frequency) and

■ **Table 1** Additional information about inputs used in evaluation: name, size n , alphabet size σ , and sizes of the corresponding synchronizing sets $|S_\tau|$ for $\tau = 256, 512, 1024,$ and 2048 .

name	n	σ	$ S_{256} $	$ S_{512} $	$ S_{1024} $	$ S_{2048} $
dblp.xml	296 135 874	97	2 304 012	1 153 799	578 703	288 758
dna	403 927 746	16	3 141 655	1 575 668	788 113	393 823
english.1024MB	1 073 741 824	239	8 354 560	4 187 042	2 095 466	1 047 924
proteins	1 184 051 855	27	9 215 429	4 616 809	2 341 374	1 155 364
sources	210 866 607	230	1 640 498	821 450	417 457	205 792
cere	461 286 644	5	31 619 034	26 205 215	20 013 847	12 699 848
coreutils	205 281 778	236	1 596 863	800 826	422 577	205 539
Escherichia_Coli	112 689 515	15	876 378	439 647	222 292	109 797
einstein.de.txt	92 758 441	117	721 365	361 616	181 080	90 702
einstein.en.txt	467 626 544	139	3 640 253	1 823 454	998 765	547 184
influenza	154 808 555	15	1 204 011	604 046	301 601	151 142
kernel	257 961 616	160	2 008 714	1 006 209	505 080	252 391

cache sizes: 32 KB L1, 256 KB L2, 25.6 MB L3 cache) and 64 GB RAM. All algorithms are sequential and only use a single core, and the results are averages of five runs. The code was compiled with GCC 9.2.0 and compiler flags `-O3`, `-ffast-math`, and `-march=native`.

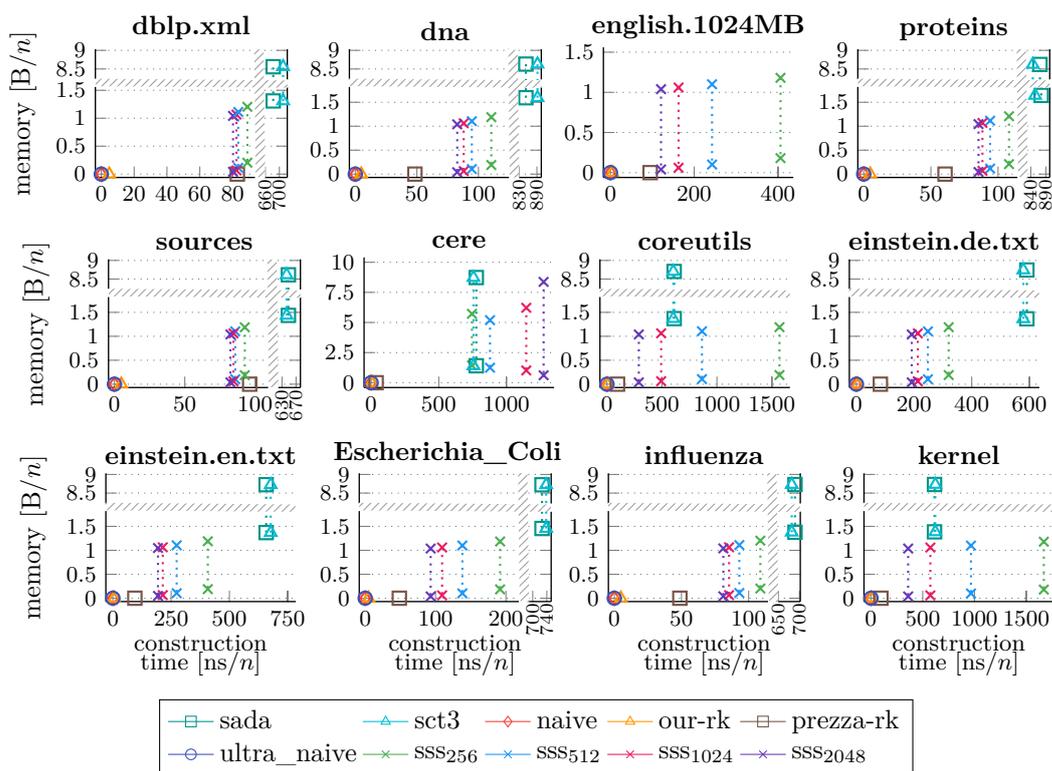
We ran all algorithms on the input shown in Tbl. 1. The texts were taken from the well-established benchmark suite *Pizza&Chili* [17]. Data sets from the top half of the table are from the *regular corpus*, whereas those from the bottom half are from the *repetitive corpus*.

Apart from several characteristics of the input texts, like length and alphabet size, Tbl. 1 also shows the sizes of the resulting string synchronizing sets for different values of τ . These numbers confirm our claim from Sect. 3.2 that the additional measures for keeping the SSS small are not necessary in practice, as it can be observed that the growth of SSS size is almost always perfectly proportional to the decrease of the values τ (growth rate only slightly less than 1). The only notable exception is the data set “cere”, one of the highly repetitive texts, which contains long runs (of the character N). There, for example, when halving τ from 512 to 256, the SSS grows only by a factor of ≈ 1.21 (instead of the expected ≈ 2).

5.1 Construction and Space Usage

Fig. 2 shows the construction times and the space usage of the data structures. For the SSS and CST data structures, that figure actually shows two numbers: the final size of the data structure (bottom mark) and the memory peak during construction time (top mark). The difference of these two numbers is therefore the additional space at construction time, which results from the intermediate steps using additional data structures (as described in Sect. 4.4 for SSS). The other data structures need no extra space during construction, as they are either in-place (“prezza-rk” and “our-rk”) or do not do any preprocessing at all (“ultra_naive” and “naive”). Also, “our-rk” is significantly faster to build than “prezza-rk” on all inputs; it is 19.76 times faster on “sources”. The CST could not be computed for “english.1024MB”. Overall, the CST data structures require the most memory. They are also the slowest to construct on all texts but “cere”, “coreutils”, and “kernel”, where SSS is slower for some τ .

The first thing to note is that both time and space requirements grow for SSS with decreasing parameter τ . For space, this is what one would expect immediately, whereas for

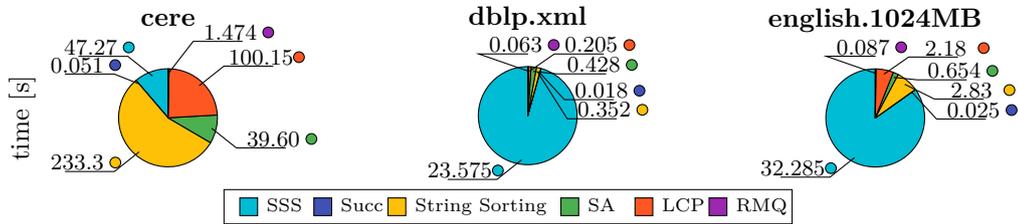


■ **Figure 2** Construction time in seconds and memory consumption in byte per character of the input for the LCE data structures. The upper mark is the memory peak during construction, and the lower mark is the memory required for the final data structure. (Needs colors for viewing.)

321 the running time, this needs further explanation: the times for *sorting* the characters of T' ,
 322 which are substrings of T , depend both on their length 3τ and their number $|S|$. One could
 323 now conjecture that in all cases, roughly the same amount of characters has to be sorted,
 324 resulting in construction times mostly independent of τ . However, as our string sorting
 325 procedure is *prefix aware* (it only considers the strings up to their distinguishing prefixes),
 326 the 3τ -long substrings are usually not inspected in full. This implies that the number of
 327 strings can have a higher impact on the running times than their total length, despite the
 328 fact that the total length of the strings remains the same (e.g. “english.1024MB”). For other
 329 texts (like “dblp.xml”), the construction times are very similar for different values of τ , but
 330 are generally faster with growing τ . The notable exception is again the data set “cere”, where
 331 the order on the time-axis is reversed (but not on the space-axis). This can be explained as
 332 follows: we already said before that the growth of the size of the SSS is much less pronounced
 333 with decreasing values of τ than for the other data sets. The string sorter is likely to inspect
 334 a number of characters proportional to the length of the strings. Therefore, the sorting times
 335 for the cere-substrings are influenced more by their total lengths than for the other data sets.

336 The construction of the SSS structure can be broken down into several steps. Fig. 3
 337 shows examples of how much time is needed for the individual components (building SSS,
 338 sorting the length- 3τ substrings, building the successor data structure, the suffix array, the
 339 LCP array, and the RMQ data structure). Some components differ significantly from text to
 340 text: for “cere”, the string sorting takes almost half of the total time, whereas for “dblp.xml”
 341 (and partly also “english”), the construction of the SSS itself takes almost all of the time.

342 The charts also indicate if and where further engineering efforts can pay off: for example,
 343 improving the successor data structure further will most likely neither speed up the final
 344 construction times significantly nor improve the space usage. On the contrary, speeding up
 345 string sorting or SSS construction (e.g., by parallelization) could result in an overall speedup.



■ **Figure 3** Snapshots of construction time for different phases of sss_{512} (same as sss_{512}^p).

346 5.2 Query Times

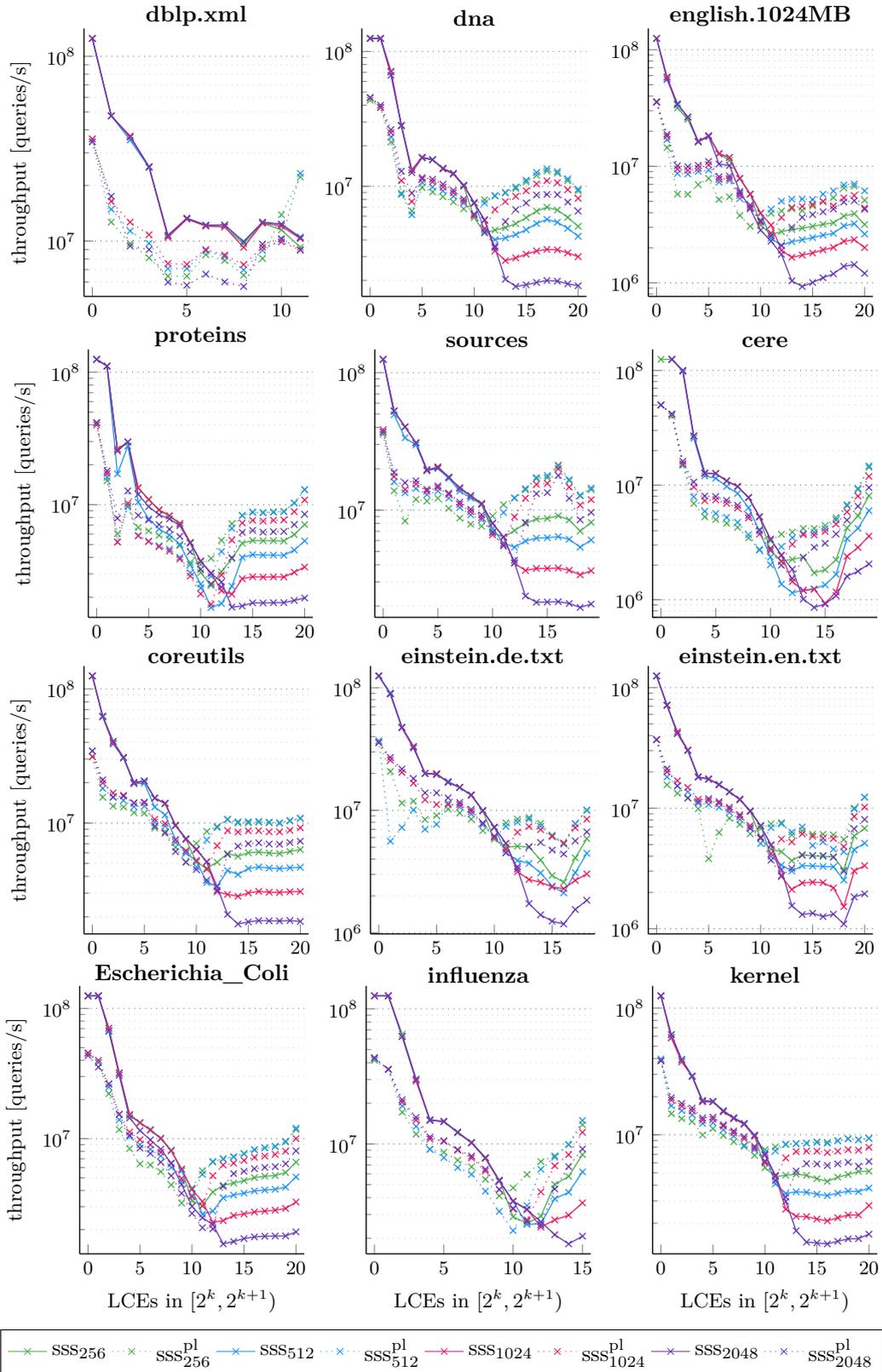
347 Now, we are interested in query times for different LCE values. To this end, we precomputed
 348 queries whose results are in the interval $[2^k, 2^{k+1})$ for all k , allowing us to get a detailed look
 349 at query times depending on the LCE value and use the same queries for all data structures.

350 We now first evaluate the query throughput of the different SSS data structures (with
 351 varying values for τ). The results are shown in Fig. 4, grouped by the lengths of the actual
 352 LCEs. It can be observed that the initial scanning of characters generally results in a visible
 353 drop of the throughput for longer LCEs; smaller values of τ lead to faster query times; the
 354 prefer-long variants are indeed faster for longer LCEs, whereas they are slower for shorter
 355 ones; and the throughput of all variants does not drop below a threshold for very long LCEs.
 356 Indeed, for many data sets, very long LCEs become faster, which can be explained by our
 357 method for answering RMQs: for very long LCEs, the corresponding LCP values are likely to
 358 be close together in the LCP array, so our query procedure is likely to use the fast scanning
 359 for the minimum instead of invoking the heavy $O(1)$ -time machinery. Given that the results
 360 for $\tau = 512$ and $\tau = 256$ are almost equally fast (in particular for “prefer-long”), we choose
 361 the larger of the two values for the following tests, as it results in a smaller data structure.

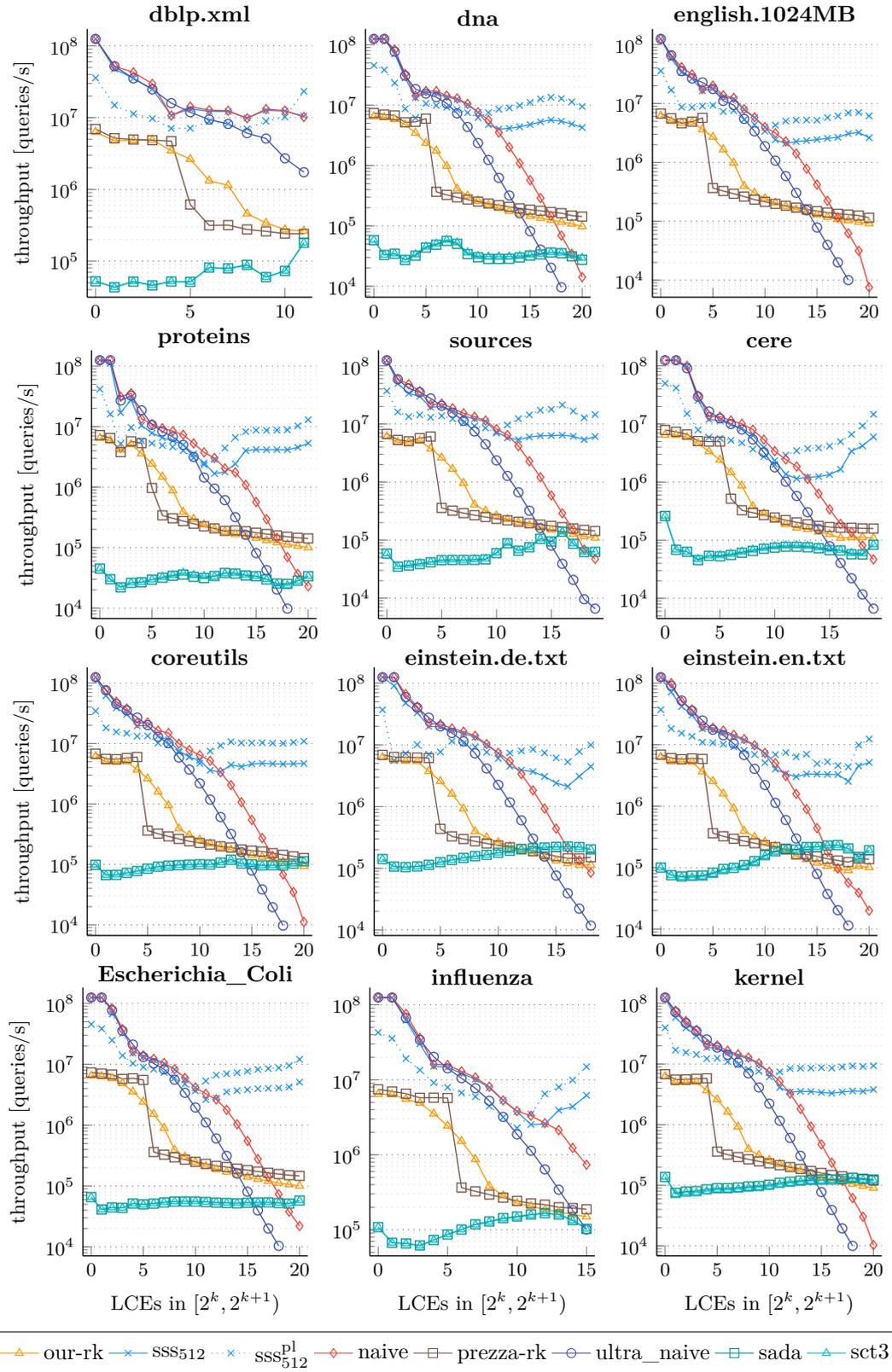
362 Fig. 5 shows the query throughput of all data structures (for SSS only with $\tau = 512$).
 363 We observe that “ultra_naive” is always slower than “naive” (this comes at no surprise);
 364 for short LCEs (roughly up to 2^8), “ultra_naive”, “naive”, and “ sss_{512} ” are fastest due to
 365 the simple initial scanning; for longer LCEs (greater than 2^{12}), “ sss_{512}^p ” becomes the fastest
 366 data structure (also faster than “ sss_{512} ”, which becomes faster than the naive approaches
 367 for LCEs longer than $\approx 2^{12}$); “our-rk” is slower than Prezza’s original implementation
 368 “prezza-rk” except for medium sized LCEs around 2^5 – 2^{10} ; and SSS’s are much faster than
 369 the fingerprinting or naive methods for LCEs longer than $\approx 2^{12}$. The CST data structures
 370 “sada” and “sct3” are of similar speed and faster than “ultra_naive” only for long LCEs,
 371 around 2^{16} on most inputs. On the repetitive texts “einstein.de.txt” and “einstein.en.txt”
 372 both are faster than the fingerprint data structures for medium-size LCEs around $\approx 2^{13}$.

373 6 Conclusions

374 We conclude from the experiments that string synchronizing sets (SSS) are the method
 375 of choice if their 10–20% overhead on top of the original text size fits into RAM, as they
 376 naturally combine the best of two worlds: they answer short LCEs equally fast as naive



■ **Figure 4** Comparing query throughput of our SSS LCE data structures.



■ **Figure 5** Query throughput of the LCE data structures, only including the fastest SSS ($\tau = 512$).

377 scanning methods, but are much faster for long LCEs and have a guaranteed worst-case
378 query time. This threshold is much earlier than previously conjectured: even for LCE values
379 larger than $\approx 2^9$, the additional effort for constructing the data structure pays off. If even the
380 little extra space for the SSS is too much and worst case query times have to be guaranteed
381 for long LCEs, then one should use an in-place fingerprinting method.

382 ——— References ———

- 383 1 Stephen Alstrup, Gerth Stølting Brodal, and Theis Rauhe. Pattern matching in dynamic
384 texts. In *11th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 819–828.
385 ACM/SIAM, 2000. URL: <http://dl.acm.org/citation.cfm?id=338219.338645>.
- 386 2 Amihood Amir, Moshe Lewenstein, and Ely Porat. Faster algorithms for string matching with
387 k mismatches. *Journal of Algorithms*, 50(2):257–275, 2004. doi:10.1016/S0196-6774(03)
388 00097-X.
- 389 3 Johannes Bahne, Nico Bertram, Marvin Böcker, Jonas Bode, Johannes Fischer, Hermann
390 Foot, Florian Grieskamp, Florian Kurpicz, Marvin Löbel, Oliver Magiera, Rosa Pink, David
391 Piper, and Christopher Poeplau. Sacabench: Benchmarking suffix array construction. In
392 *26th International Symposium on String Processing and Information Retrieval (SPIRE)*,
393 volume 11811 of *Lecture Notes in Computer Science*, pages 407–416. Springer, 2019. doi:
394 10.1007/978-3-030-32686-9_29.
- 395 4 Hideo Bannai, Tomohiro I, Shunsuke Inenaga, Yuto Nakashima, Masayuki Takeda, and
396 Kazuya Tsuruta. The "runs" theorem. *SIAM J. Comput.*, 46(5):1501–1514, 2017. doi:
397 10.1137/15M1011032.
- 398 5 Hideo Bannai, Shunsuke Inenaga, and Dominik Köppl. Computing all distinct squares in
399 linear time for integer alphabets. In *28th Annual Symposium on Combinatorial Pattern
400 Matching (CPM)*, volume 78 of *LIPICs*, pages 22:1–22:18. Schloss Dagstuhl - Leibniz-Zentrum
401 für Informatik, 2017. doi:10.4230/LIPICs.CPM.2017.22.
- 402 6 Philip Bille, Anders Roy Christiansen, Patrick Hagge Cording, and Inge Li Gørtz. Finger
403 search in grammar-compressed strings. *Theory of Computing Systems*, 62(8):1715–1735, 2018.
404 doi:10.1007/s00224-017-9839-9.
- 405 7 Philip Bille, Johannes Fischer, Inge Li Gørtz, Tsvi Kopelowitz, Benjamin Sach, and
406 Hjalte Wedel Vildhøj. Sparse text indexing in small space. *ACM Transactions on Algo-
407 rithms*, 12(3):39:1–39:19, 2016. doi:10.1145/2836166.
- 408 8 Philip Bille, Inge Li Gørtz, Patrick Hagge Cording, Benjamin Sach, Hjalte Wedel Vildhøj, and
409 Søren Vind. Fingerprints in compressed strings. *Journal of Computer and System Sciences*,
410 86:171–180, 2017. doi:10.1016/j.jcss.2017.01.002.
- 411 9 Philip Bille, Inge Li Gørtz, Mathias Bæk Tejs Knudsen, Moshe Lewenstein, and Hjalte Wedel
412 Vildhøj. Longest common extensions in sublinear space. In *26th Annual Symposium on
413 Combinatorial Pattern Matching (CPM)*, volume 9133 of *Lecture Notes in Computer Science*,
414 pages 65–76. Springer, 2015. doi:10.1007/978-3-319-19929-0_6.
- 415 10 Philip Bille, Inge Li Gørtz, Benjamin Sach, and Hjalte Wedel Vildhøj. Time-space trade-
416 offs for longest common extensions. *Journal of Discrete Algorithms*, 25:42–50, 2014. doi:
417 10.1016/j.jda.2013.06.003.
- 418 11 Timo Bingmann, Andreas Eberle, and Peter Sanders. Engineering parallel string sorting.
419 *Algorithmica*, 77(1):235–286, 2017. doi:10.1007/s00453-015-0071-1.
- 420 12 Or Birenzweige, Shay Golan, and Ely Porat. Locally consistent parsing for text indexing in
421 small space. In *30th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages
422 607–626. SIAM, 2020. doi:10.1137/1.9781611975994.37.
- 423 13 Maxime Crochemore, Roman Kolpakov, and Gregory Kucherov. Optimal bounds for computing
424 α -gapped repeats. *Inf. Comput.*, 268, 2019. doi:10.1016/j.ic.2019.104434.
- 425 14 Roman Dementiev, Lutz Kettner, Jens Mehnert, and Peter Sanders. Engineering a sorted list
426 data structure for 32 bit key. In *Sixth Workshop on Algorithm Engineering and Experiments*

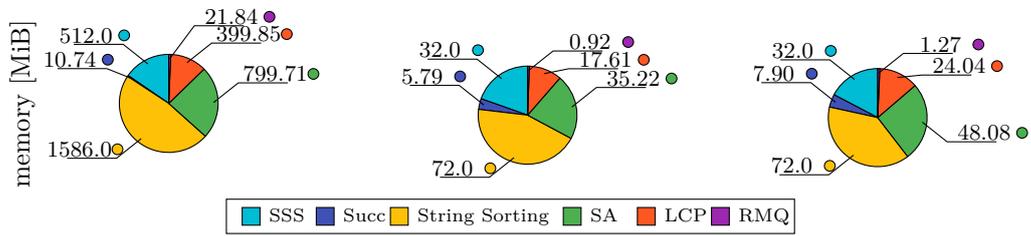
- 427 (ALLENEX) and the First Workshop on Analytic Algorithmics and Combinatorics (ANALCO),
428 pages 142–151. SIAM, 2004.
- 429 **15** Martin Farach-Colton, Paolo Ferragina, and S. Muthukrishnan. On the sorting-complexity
430 of suffix tree construction. *Journal of the ACM*, 47(6):987–1011, 2000. doi:10.1145/355541.
431 355547.
- 432 **16** Héctor Ferrada and Gonzalo Navarro. Improved range minimum queries. *J. Discrete Algorithms*,
433 43:72–80, 2017. doi:10.1016/j.jda.2016.09.002.
- 434 **17** Paolo Ferragina and Gonzalo Navarro. Pizza&Chili corpus: Compressed indexes and their
435 testbeds. URL: <http://pizzachili.dcc.uchile.cl/>.
- 436 **18** Johannes Fischer and Volker Heun. Theoretical and practical improvements on the rmq-
437 problem, with applications to LCA and LCE. In *17th Annual Symposium on Combinatorial
438 Pattern Matching (CPM)*, volume 4009 of *Lecture Notes in Computer Science*, pages 36–48.
439 Springer, 2006. doi:10.1007/11780441_5.
- 440 **19** Johannes Fischer, Tomohiro I, and Dominik Köppl. Deterministic sparse suffix sorting on
441 rewritable texts. In *12th Latin American Symposium on Theoretical Informatics (LATIN)*,
442 volume 9644 of *Lecture Notes in Computer Science*, pages 483–496. Springer, 2016. doi:
443 10.1007/978-3-662-49529-2_36.
- 444 **20** Johannes Fischer, Veli Mäkinen, and Gonzalo Navarro. Faster entropy-bounded compressed
445 suffix trees. *Theor. Comput. Sci.*, 410(51):5354–5364, 2009. doi:10.1016/j.tcs.2009.09.012.
- 446 **21** Johannes Fischer, Veli Mäkinen, and Niko Välimäki. Space efficient string mining under
447 frequency constraints. In *8th IEEE International Conference on Data Mining (ICDM)*, pages
448 193–202. IEEE Computer Society, 2008. doi:10.1109/ICDM.2008.32.
- 449 **22** Zvi Galil and Raffaele Giancarlo. Improved string matching with k mismatches. *SIGACT
450 News*, 17(4):52–54, 1986. doi:10.1145/8307.8309.
- 451 **23** Paweł Gawrychowski, Adam Karczmarz, Tomasz Kociumaka, Jakub Łacki, and Piotr
452 Sankowski. Optimal dynamic strings. In *29th Annual ACM-SIAM Symposium on Dis-
453 crete Algorithms (SODA)*, pages 1509–1528. SIAM, 2018. arXiv:1511.02612v2, doi:10.1137/
454 1.9781611975031.99.
- 455 **24** Paweł Gawrychowski and Tomasz Kociumaka. Sparse suffix tree construction in optimal time
456 and space. In *28th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages
457 425–439. SIAM, 2017. doi:10.1137/1.9781611974782.27.
- 458 **25** Simon Gog, Timo Beller, Alistair Moffat, and Matthias Petri. From theory to practice: Plug
459 and play with succinct data structures. In Joachim Gudmundsson and Jyrki Katajainen,
460 editors, *13th International Symposium on Experimental Algorithms, SEA 2014*, volume 8504
461 of *LNCS*, pages 326–337. Springer, 2014. doi:10.1007/978-3-319-07959-2_28.
- 462 **26** Dan Gusfield. *Algorithms on Strings, Trees, and Sequences*. Cambridge UP, 1997. doi:
463 10.1017/cbo9780511574931.
- 464 **27** Tomohiro I. Longest common extensions with recompression. In *28th Annual Symposium
465 on Combinatorial Pattern Matching (CPM)*, volume 78 of *LIPICs*, pages 18:1–18:15. Schloss
466 Dagstuhl–Leibniz-Zentrum für Informatik, 2017. doi:10.4230/LIPICs.CPM.2017.18.
- 467 **28** Tomohiro I, Wataru Matsubara, Kouji Shimohira, Shunsuke Inenaga, Hideo Bannai, Masayuki
468 Takeda, Kazuyuki Narisawa, and Ayumi Shinohara. Detecting regularities on grammar-
469 compressed strings. *Information and Computation*, 240:74–89, 2015. doi:10.1016/j.ic.2014.
470 09.009.
- 471 **29** Lucian Ilie and Liviu Tinta. Practical algorithms for the longest common extension problem.
472 In *16th International Symposium on String Processing and Information Retrieval (SPIRE)*,
473 volume 5721 of *Lecture Notes in Computer Science*, pages 302–309. Springer, 2009. doi:
474 10.1007/978-3-642-03784-9_30.
- 475 **30** Juha Kärkkäinen, Peter Sanders, and Stefan Burkhardt. Linear work suffix array construction.
476 *Journal of the ACM*, 53(6):918–936, 2006. doi:10.1145/1217856.1217858.
- 477 **31** Richard M. Karp and Michael O. Rabin. Efficient randomized pattern-matching algorithms.
478 *IBM Journal of Research and Development*, 31(2):249–260, 1987. doi:10.1147/rd.312.0249.

- 479 32 Dominik Kempa and Tomasz Kociumaka. String synchronizing sets: sublinear-time BWT
480 construction and optimal LCE data structure. In *51st Annual ACM SIGACT Symposium on*
481 *Theory of Computing (STOC)*, pages 756–767. ACM, 2019. doi:10.1145/3313276.3316368.
- 482 33 Tomasz Kociumaka, Jakub Radoszewski, Wojciech Rytter, and Tomasz Waleń. Internal pattern
483 matching queries in a text and applications. In *26th Annual ACM-SIAM Symposium on Discrete*
484 *Algorithms (SODA)*, pages 532–551. SIAM, 2015. doi:10.1137/1.9781611973730.36.
- 485 34 Dmitry Kosolobov. Tight lower bounds for the longest common extension problem. *Inf.*
486 *Process. Lett.*, 125:26–29, 2017. doi:10.1016/j.ipl.2017.05.003.
- 487 35 Gad M. Landau, Eugene W. Myers, and Jeanette P. Schmidt. Incremental string comparison.
488 *SIAM J. Comput.*, 27(2):557–582, 1998. doi:10.1137/S0097539794264810.
- 489 36 Gad M. Landau and Uzi Vishkin. Efficient string matching with k mismatches. *Theor. Comput.*
490 *Sci.*, 43:239–249, 1986. doi:10.1016/0304-3975(86)90178-7.
- 491 37 Gad M. Landau and Uzi Vishkin. Fast string matching with k differences. *Journal of Computer*
492 *and System Sciences*, 37(1):63–78, 1988. doi:10.1016/0022-0000(88)90045-1.
- 493 38 Yoshiaki Matsuoka, Tomohiro I, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda.
494 Semi-dynamic compact index for short patterns and succinct van Emde Boas tree. In *26th*
495 *Annual Symposium on Combinatorial Pattern Matching (CPM)*, volume 9133 of *Lecture Notes*
496 *in Computer Science*, pages 355–366. Springer, 2015. doi:10.1007/978-3-319-19929-0_30.
- 497 39 Kurt Mehlhorn, R. Sundar, and Christian Uhrig. Maintaining dynamic sequences under equality
498 tests in polylogarithmic time. *Algorithmica*, 17(2):183–198, 1997. doi:10.1007/BF02522825.
- 499 40 Yuta Mori. sais: An implementation of the induced sorting algorithm. URL: <https://sites.google.com/site/yuta256/>.
- 501 41 Gonzalo Navarro. *Compact Data Structures - A Practical Approach*. Cambridge University
502 Press, 2016. doi:10.1017/cbo9781316588284.
- 503 42 Takaaki Nishimoto, Tomohiro I, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda.
504 Fully dynamic data structure for LCE queries in compressed space. In *41st International*
505 *Symposium on Mathematical Foundations of Computer Science (MFCS)*, volume 58 of *LIPICs*,
506 pages 72:1–72:15. Schloss Dagstuhl–Leibniz-Zentrum für Informatik, 2016. doi:10.4230/
507 LIPICs.MFCS.2016.72.
- 508 43 Cinzia Pizzi. Missmax: alignment-free sequence comparison with mismatches through filtering
509 and heuristics. *Algorithms Mol. Biol.*, 11:6, 2016. doi:10.1186/s13015-016-0072-x.
- 510 44 Nicola Prezza. In-place sparse suffix sorting. In *29th Annual ACM-SIAM Symposium on*
511 *Discrete Algorithms (SODA)*, pages 1496–1508. SIAM, 2018. doi:10.1137/1.9781611975031.
512 98.
- 513 45 Wei Quan, Bo Liu, and Yadong Wang. SALT: a fast, memory-efficient and snp-aware short
514 read alignment tool. In *IEEE International Conference on Bioinformatics and Biomedicine*
515 *(BIBM)*, pages 1774–1779. IEEE, 2019. doi:10.1109/BIBM47256.2019.8983162.
- 516 46 Süleyman Cenk Sahinalp and Uzi Vishkin. On a parallel-algorithms method for string matching
517 problems. In *Second Italian Conference on Algorithms and Complexity (CIAC)*, volume 778 of
518 *LNCS*, pages 22–32. Springer, 1994. doi:10.1007/3-540-57811-0_3.
- 519 47 Avi Srivastava, Hirak Sarkar, Nitish Gupta, and Robert Patro. Rapmap: a rapid, sensitive
520 and accurate tool for mapping rna-seq reads to transcriptomes. *Bioinform.*, 32(12):192–200,
521 2016. doi:10.1093/bioinformatics/btw277.
- 522 48 Yuka Tanimura, Tomohiro I, Hideo Bannai, Shunsuke Inenaga, Simon J. Puglisi, and Masayuki
523 Takeda. Deterministic sub-linear space LCE data structures with efficient construction. In
524 *27th Annual Symposium on Combinatorial Pattern Matching (CPM)*, volume 54 of *LIPICs*,
525 pages 1:1–1:10. Schloss Dagstuhl–Leibniz-Zentrum für Informatik, 2016. doi:10.4230/LIPICs.
526 CPM.2016.1.
- 527 49 Peter van Emde Boas. Preserving order in a forest in less than logarithmic time and linear
528 space. *Inf. Process. Lett.*, 6(3):80–82, 1977. doi:10.1016/0020-0190(77)90031-X.

16 Practical performance of data structures for LCE queries

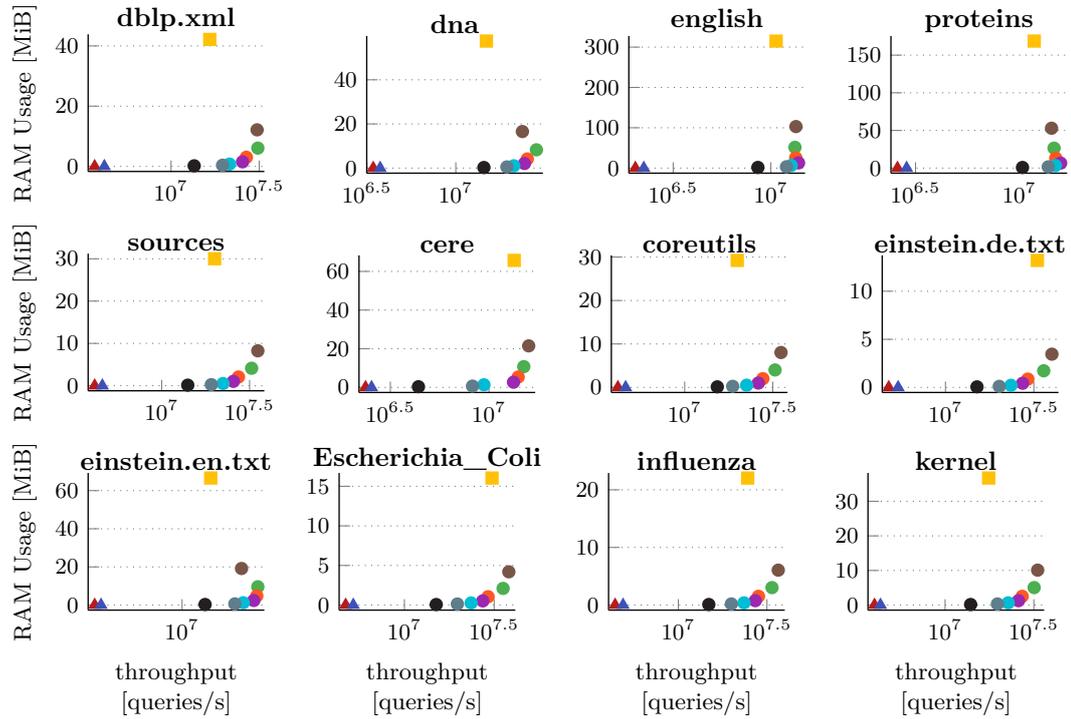
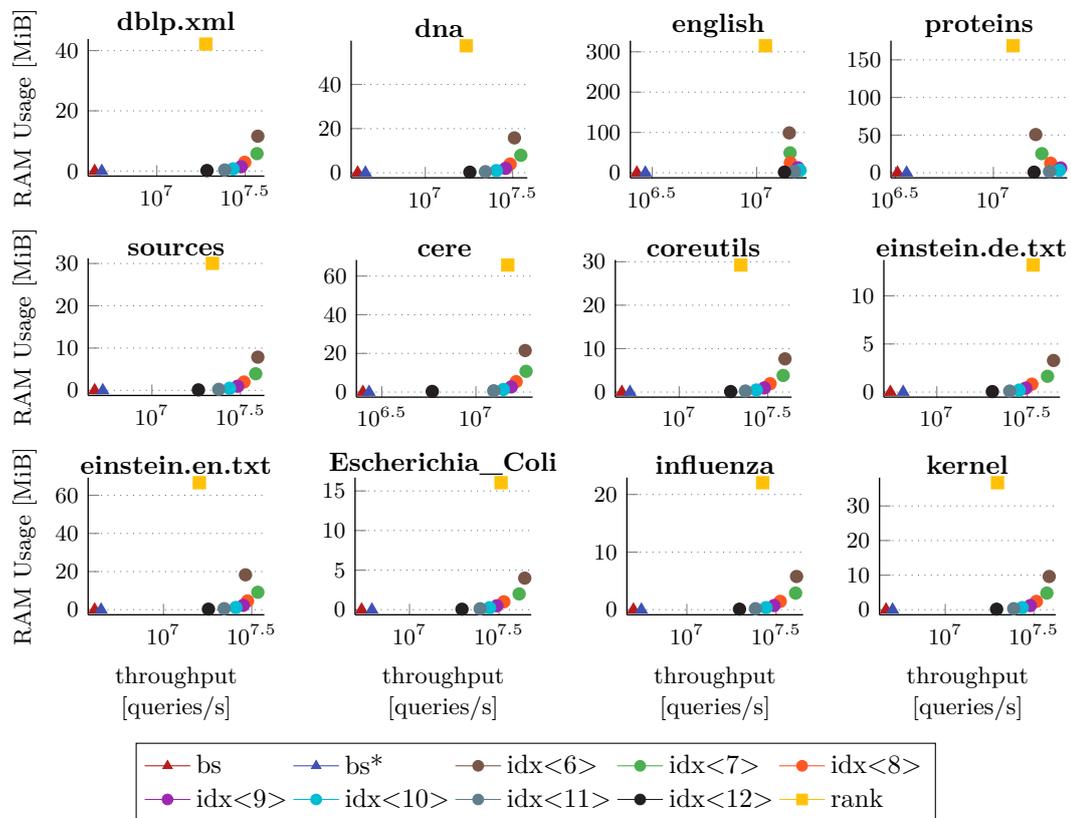
- 529 50 Chen Zhou, Hao Chi, Leheng Wang, You Li, Yan-Jie Wu, Yan Fu, Ruixiang Sun, and Si-Min
530 He. Speeding up tandem mass spectrometry-based database searching by longest common
531 prefix. *BMC Bioinform.*, 11:577, 2010. doi:10.1186/1471-2105-11-577.

532 **A Detailed Memory During Construction of SSS LCE Data Structure**



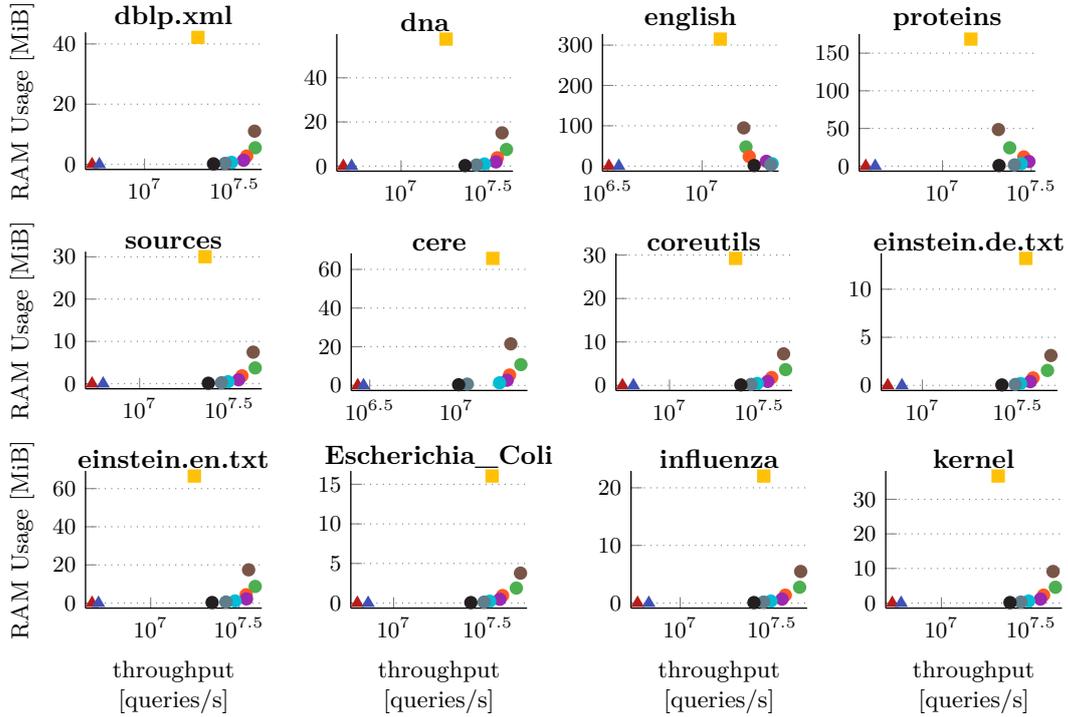
■ **Figure 6** Snapshots of memory during construction for different phases of sss_{512} (same as sss_{512}^{p1}).

533 **B Detailed Results for Predecessor Queries**

(a) $\tau = 256$ (b) $\tau = 512$ 

■ **Figure 7** Query throughput versus space usage for static successor data structures on SSS's.

(a) $\tau = 1024$



(b) $\tau = 2048$

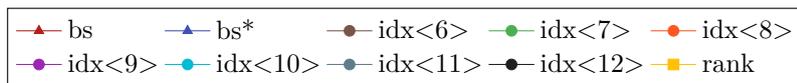
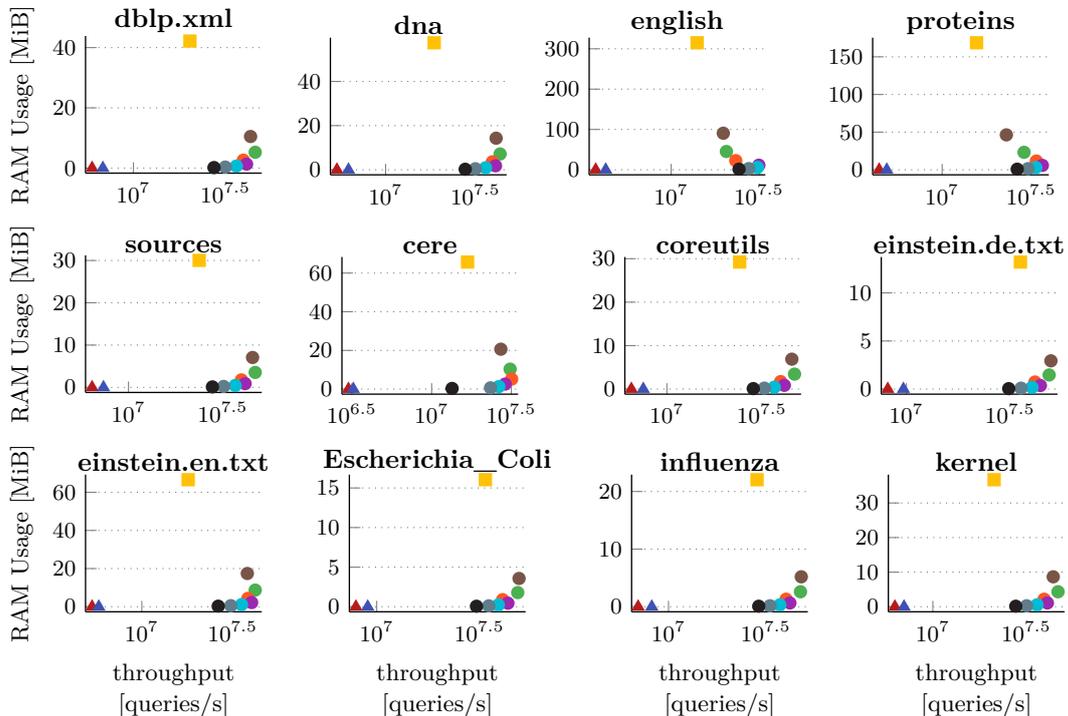


Figure 8 Query throughput versus space usage for static successor data structures on SSS's (cont.).