

Improved Characters Distance Sampling for Online and Offline Text Searching¹

Simone Faro^a, Francesco Pio Marino^a, Arianna Pavone^b

^a*Department of Mathematics and Computer Science, University of Catania, Viale A.Doria n.6, Catania, 95125, Italy*

^b*Department of Cognitive Science, University of Messina, Via Concezione n.6, Messina, 98121, Italy*

Abstract

Sampled string matching is a very effective technique to reduce the search time for a pattern within a text at the cost of a small amount of additional memory, used for storing a partial index of the text. This approach has recently received some interest and has been applied to improve both online and offline string matching solutions, improving standard solutions by more than 50%. However, this improvement is currently only achievable in the case of texts on large-sized alphabets, and remains small (or absent) in the case of small-sized alphabets. In this article we propose an extension of the approach to text-sampling, known as Character Distance Sampling, to the case of small alphabets, obtaining an improvement of up to 98% compared to standard solutions in the case of online string matching. We also extend this approach to the case of offline string matching, introducing a sampled version of the suffix array, obtaining performances up to 5 times higher than the search obtained on the standard suffix array. Differently from what has been done by previous solutions, our idea is not based on the reduction of the number of indexed suffixes, but on the construction of the index directly on the sampled text.

Keywords: Text processing, experimental algorithms, string matching

¹This paper is based on preliminary results appeared in Proceedings of the 22nd Italian Conference on Theoretical Computer Science (ICTCS 2021) [13] and in Proceedings of the 25th Prague Stringology Conference (PSC 2021) [11]

1. Introduction

Exact string matching is a fundamental problem in computer science and in the wide domain of text processing. It consists in finding all the occurrences of a given pattern x in a large text y where characters of both sequences are drawn from an alphabet Σ .

It is a fundamental problem in computer science with applications in many other fields, like natural language processing and information retrieval. It is also a critical problem in computational molecular biology and plays a very important role in biological sequences analysis, mainly due to the constantly growing amount of molecular data extracted from living organisms. For this reason sequence matching techniques play a very important role in various applications in computational biology for data analysis. In addition, as the size of data increases, the space required to store this data and the data structures useful for solving the problem is also constantly increasing, which is why it is necessary to adopt new efficient approaches that can drastically reduce the space used while preserving the effectiveness of the search.

Applications require two kinds of solutions: *online* and *offline* string matching. Solutions based on the first approach assume that the text is not preprocessed and thus they need to scan the input sequence *online*, when searching. Differently, solutions based on the second approach try to drastically speed up searching by preprocessing the text and building a data structure that allows searching in time proportional to the length of the pattern. For this reason such kind of problem is known as *indexed searching*. *Sampled string matching* is a technique that has recently received interest and which is halfway between the two just described solutions, allowing both to be improved. Its goal is to significantly cut down the space requirements of indexed matching, on the one hand, and drastically reduce searching time for the online solutions, on the other hand. However sampled string matching allows so far to search efficiently only in the case of natural language texts or, in general, when searching on input sequences over large alphabets, while its performances degrade when the size of the underlying alphabets decreases.

In this paper, we present an extension of a previous approach [12], called Character Distance Sampling (CDS), in two separate directions.

- We extend CDS to small alphabets obtaining a more feasible solution in the case of biological data like genome or protein sequences and, in general, in the case of small alphabets. Our proposed approach

37 makes use of condensed characters in order to enlarge the size of the
38 underlying alphabet and, as a result, speed up the searching process
39 and reduce the space consumption of the resulting sampled text.

- 40 • We adapt the proposed sampling approach also in the case of the offline
41 search by developing a sampled variant of the suffix array based on
42 the CDS approach. Differently from what has been done by previous
43 solutions, our idea is not based on the reduction of the number of
44 indexed suffixes, but on the construction of the index directly on the
45 sampled text.

46 From our experimental results it turns out that the use of condensed
47 alphabets leads to reduce the space consumption up to 80% and to speed
48 up the online searching process up to 98%, significantly improving the re-
49 sults obtained by the previous text sampling approach. We also conducted
50 experimental tests for offline search, obtaining, also in this case, significant
51 improvements. Specifically, the new approach presented in this paper allows
52 a reduction of search times up to 5 times, compared to the times obtained
53 using the standard data structure.

54 The paper is organized as follows. First we briefly review previous so-
55 lutions related with our work making appropriate references to the most
56 relevant literature. In Section 3 we briefly review the Characters Distance
57 Sampling approach introduced by Faro *et al.* [12] and extend it to condensed
58 alphabets. Then, in Section 4, we show how to apply our sampling approach
59 to the online string matching case, while its application to the case of offline
60 string matching is presented in Section 5. Finally, in Section 6, we present
61 experimental results and draw our conclusions in Section 7.

62 2. Related Results

63 Formally, the *exact string matching* problem consists in finding all the
64 occurrences of a given pattern x , of length m , in a large text y , of length n ,
65 where characters of both sequences are drawn from an alphabet Σ of size σ .

66 *Online string matching* solutions assume that the text is not preprocessed
67 and thus they need to scan the input sequence *online*, when searching. Their
68 worst case time complexity is $\Theta(n)$, and was achieved for the first time by
69 the well known Knuth-Morris-Pratt (KMP) algorithm [24], while the optimal
70 average time complexity of the problem is $\Theta(n \log_{\sigma}(m)/m)$ [38], achieved for
71 example by the Backward-Dawg-Matching (BDM) algorithm [6].

72 Many string matching solutions have been also developed in order to
73 obtain sub-linear performance in practice [9]. Among them the Boyer-Moore-
74 Horspool (BMH) algorithm [2, 20] deserves a special mention, since it has
75 inspired much work. Memory requirements of this class of algorithms are
76 very low and generally limited to a precomputed table of size $O(m\sigma)$ or
77 $O(\sigma^2)$ [9]. However their searching time is always proportional to the length
78 of the text and thus their performances may stay poor in many practical
79 cases, especially for huge texts and short patterns.²

80 Differently, solutions based on *indexed searching* try to drastically speed
81 up searching by preprocessing the text and building a data structure that
82 allows searching in time proportional to the length of the pattern. The
83 literature in this research area is truly extensive and citing all the solutions
84 proposed over the years would go beyond the scope of this paper. However,
85 among the most efficient solutions to such problem we mention those based on
86 suffix trees [9], which find all occurrences in $O(m+occ)$ -worst case time, those
87 based on suffix arrays [2], which solve the problem in $O(m + \log n + occ)$ [2],
88 where occ is the number of occurrences of x in y , and those based on the FM-
89 index [15] (Full-text index in Minute space), which is a compressed full-text
90 substring index based on the Burrows-Wheeler Transform (BWT) allowing
91 compression of the input text while still permitting fast substring queries.
92 However, despite their optimal time performance³, space requirements of
93 full-index data structures, as suffix-trees and suffix-arrays, are from 4 to 20
94 times the size of the text.

95 While the size of a compressed indexes, as the FM-Index [15], is typically
96 less than the size of the text, it turns out that their space requirement is too
97 large for many practical applications.

98 An alternative solution to full indexes is to compress the input text and
99 search online directly the compressed data in order to speed-up the search-
100 ing process using reduced extra space. Such problem, known in literature
101 as *compressed string matching*, has been widely investigated in the last few
102 years. Although efficient solutions exist for searching on standard compres-
103 sions schemes, as Ziv-Lempel [33] and Huffman [3], the best practical be-

²Search speed of an online string matching algorithm may depend on the length of the pattern. Typical search speed of a fast solution, on a modern laptop computer, goes from 1 GB/s (in the case of short patterns) to 5 GB/s (in the case of very long patterns) [4].

³Search speed of a fast offline solution do not depend on the length of the text and is typically under 1 millisecond per query.

104 haviour are achieved by ad-hoc schemes designed for allowing fast searching
105 [29, 7, 23, 36, 16]. These latter solutions use less than 70% of text size ex-
106 tra space (achieving a compression rate over 30%) and are twice as fast in
107 searching as standard online string matching algorithms. A drawback of such
108 solutions is that most of them still require significant implementation efforts
109 and a high time for each reported occurrence.

110 One of the most interesting solutions to the problem are *compact data*
111 *structures*. Such structures are equipped with native tools for handling text
112 directly in its compressed form [31]. In general, however, they are not able
113 to compress text by orders of magnitude, offering only complex functionality
114 in the space required by raw data.

115 When working on repetitive texts the BWT, featuring long runs of equal
116 consecutive symbols, has enormous potential in terms of compression [32]
117 and compact data structures benefit from this feature. These include, for
118 example, a complete index for pan-genomics read alignment using prefix-
119 free parsing [26]. A relevant compressibility measure for a repetitive text
120 is indeed the number r of runs in their BWT. Based on this measure, the
121 Run-Length FM-index [28] is able to efficiently count the number of occur-
122 rences of a pattern using $O(r)$ space and in log-logarithmic time per pattern
123 symbol. Although it can be also extended [17, 18] in order to be able to
124 locate the positions of such occurrences without using additional space, such
125 data structures are designed to be efficient only in the case of very repetitive
126 texts.

127 2.1. *Sampled String Matching*

128 An alternative solution to the problem is *sampled string matching*, in-
129 troduced in 1991 by Vishkin [37], which consists in the construction of a
130 succinct sampled version of the text (which must be maintained together
131 with the original text) and in the application of any online string matching
132 algorithm directly on the sampled sequence.

133 Although any candidate occurrence of the pattern may be found more
134 efficiently, the drawback of this approach is that any occurrence reported in
135 the sampled-text requires to be verified in the original text. Apart from this
136 point a sampled-text approach may have a lot of good features: it may be easy
137 to implement if compared with other succinct matching approaches, it may
138 require very small extra space and may allow fast searching. Additionally it
139 may also allow fast updates of the data structure.

140 Apart the theoretical result of Vishkin, the first practical solution to
141 sampled string matching has been introduced by Claude *et al.* [5] and is
142 based on an alphabet reduction. In this paper we refer to this algorithm
143 as Occurrence Text Sampling (OTS). Specifically, if we let y be the input
144 text, of length n , and let x be the input pattern, of length m , both over an
145 alphabet Σ of size σ , the main idea of the OTS approach is to select a subset
146 of the alphabet, $\hat{\Sigma} \subset \Sigma$ (the sampled alphabet), and then to construct a
147 partial-index as the subsequence of the text (the sampled text) \hat{y} , of length
148 \hat{n} , containing all (and only) the characters of the sampled alphabet $\hat{\Sigma}$. More
149 formally $\hat{y}[i] \in \hat{\Sigma}$, for all $1 \leq i \leq \hat{n}$. However, since \hat{y} contains partial
150 information, a table ρ is maintained in order to map, at regular intervals,
151 positions of the sampled text to their corresponding positions in y .

152 It turns out that the OTS approach leads to solutions which are to be up
153 to 5 times faster than standard online string matching and 2 times faster than
154 standard offline string matching on English texts. Such results are obtained
155 with an extra space requirement which is only 14% of text size.⁴

156 More recently Faro *et al.* presented a more effective sampling approach
157 based on *character distance sampling* (CDS) [12, 11], obtaining in practice
158 a speed up by a factor of up to 9 on English texts, using limited additional
159 space whose amount goes from 11% to 2.8% of the text size, with a gain in
160 searching time up to 50% if compared against the OTS approach. We will
161 describe in more detail the ideas on which the CDS approach is based in the
162 next section, in which we will extend their application to search for texts
163 on smaller alphabets using a technique based on condensed alphabets. In
164 fact, it should be emphasized that both OTS and CDS approaches to exact
165 string matching prove to work efficiently only in the case of natural language
166 texts or, in general, when searching on input sequences over large alphabets,
167 while their performances degrade when the size of the underlying alphabets
168 decreases.

⁴We also notice that some partial improvements have been thereafter presented by Grabowsky and Raniszewski [19]. They proposed a more convenient indexing suffix sampling approach, with only a minimum pattern length as a requirement. Their experiments show that the resulting solution achieves competitive time-space tradeoffs on most standard benchmark data.

169 **3. Character Distance Sampling and Condensed Alphabets**

170 In this section we briefly present the sampling approach known as Char-
 171 acter Distance Sampling (CDS) and extend it to the case of condensed al-
 172 phabets.

173 Let y be the input text, of length n , and let x be the input pattern, of
 174 length m , both over an alphabet Σ of size σ . We assume that all strings can
 175 be treated as vectors starting at position 1. Thus we refer to $x[i]$ as the i -th
 176 character of the string x , for $1 \leq i \leq m$, where m is the size of x .

177 We elect a set $C \subseteq \Sigma$ to be the *set of pivot characters*. Given this
 178 set of pivot characters we sample the text y by taking into account the
 179 distances between consecutive positions of any pivot characters $c \in C$ in y .
 180 More formally our sampling approach is based on the following definition of
 181 *position sampling* of a text.

182 **Definition 1** (Position Sampling). *Let y be a text of length n , let $C \subseteq \Sigma$
 183 be the set of pivot characters and let n_C be the number of occurrences of any
 184 $c \in C$ in the input text y .*

*First we define the position function, $\delta : \{1, \dots, n_C\} \rightarrow \{1, \dots, n\}$, where $\delta(i)$
 is the position of the i -th occurrence of any character of C in y . Formally
 we have*

- (i) $1 \leq \delta(i) < \delta(i+1) \leq n$ for each $1 \leq i \leq n_C - 1$
- (ii) $y[\delta(i)] \in C$ for each $1 \leq i \leq n_C$
- (iii) $y[\delta(i) + 1.. \delta(i+1) - 1]$ contains no $c \in C$ for each $0 \leq i \leq n_C$

185 where in (iii) we assume that $\delta(0) = 0$ and $\delta(n_C + 1) = n + 1$.

186 Then the position sampled version of y , indicated by \dot{y} , is a numeric
 187 sequence, of length n_C , defined as

$$\dot{y} = \langle \delta(1), \delta(2), \dots, \delta(n_C) \rangle. \tag{1}$$

188 **Example 1.** *Suppose $y = \text{“agaacgcagtata”}$ is a DNA sequence of length 13,
 189 over the alphabet $\Sigma = \{a, c, g, t\}$. Let $C = \{a\}$ be the set of pivot characters.
 190 Thus the position sampled version of y is $\dot{y} = \langle 1, 3, 4, 8, 11, 13 \rangle$. Specifically
 191 the first occurrence of character $a \in C$ is at position 1 ($y[1] = a$), its second
 192 occurrence is at position 3 ($y[3] = a$), and so on.*

193 **Example 2.** *As in Example 1, assume $y = \text{“agaacgcagtata”}$ is a DNA se-
 194 quence of length 13, over the alphabet $\Sigma = \{a, c, g, t\}$. Let $C = \{a, c\}$ be the*

195 set of pivot characters (now C contains two characters). Thus the position
 196 sampled version of y is $\dot{y} = \langle 1, 3, 4, 5, 7, 8, 11, 13 \rangle$. Note that in this example
 197 we simply added the two positions, 5 and 7, where the character c occurs.

198 **Definition 2** (Characters Distance Sampling). Let $C \subseteq \Sigma$ be the set of pivot
 199 characters, let $n_C \leq n$ be the number of occurrences of any pivot character
 200 in the text y and let δ be the position function of y . We define the characters
 201 distance function $\Delta(i) = \delta(i+1) - \delta(i)$, for $1 \leq i \leq n_C - 1$, as the distance
 202 between two consecutive occurrences of any pivot character in y .

203 Then the characters-distance sampled version of the text y is a numeric
 204 sequence, indicated by \bar{y} , of length $n_C - 1$ defined as

$$\begin{aligned} \bar{y} &= \langle \Delta(1), \Delta(2), \dots, \Delta(n_C - 1) \rangle \\ &= \langle \delta(2) - \delta(1), \dots, \delta(n_C) - \delta(n_C - 1) \rangle \end{aligned} \quad (2)$$

Plainly we have

$$\sum_{i=1}^{n_C-1} \Delta(i) \leq n - 1.$$

205 **Example 3.** Let $y = \text{“agaacgcagtata”}$ be a text of length 13, over the alphabet
 206 $\Sigma = \{a, c, g, t\}$. Let $C = \{a\}$ be the set of pivot characters. Thus the character
 207 distance sampling version of y is $\bar{y} = \langle 2, 1, 4, 3, 2 \rangle$. Specifically $\bar{y}[1] = \Delta(1) =$
 208 $\delta(2) - \delta(1) = 3 - 1 = 2$, while $\bar{y}[3] = \Delta(3) = \delta(4) - \delta(3) = 8 - 4 = 4$, and so
 209 on.

210 **Definition 3** (Rank of a character). Let x be a pattern of length m , and
 211 let $c \in \Sigma$. We define $\phi : \Sigma \rightarrow \{0..m\}$ as the function which associates any
 212 character of the text with the number of its occurrences in x . The rank of
 213 the character c is the position of c in the alphabet Σ , if we assume that all
 214 characters are sorted by their $\phi(c)$ values in non increasing order. If two
 215 characters of the alphabet have the same number of occurrences, then their
 216 relative order is irrelevant. To avoid confusion, we assume that they are
 217 arranged in lexicographic order. More formally the rank of c is given by the
 218 cardinality of the set $\{k \in \Sigma \mid \phi(k) > \phi(c) \text{ or } (\phi(k) = \phi(c) \text{ and } k > c)\} + 1$

219 **Example 4.** Let again $y = \text{“agaacgcagtata”}$ be a text of length 13, over the
 220 alphabet $\Sigma = \{a, c, g, t\}$. The values associated by the function ϕ to the four
 221 characters of the alphabet are $\phi(a) = 6$, $\phi(c) = 2$, $\phi(g) = 3$ and $\phi(t) = 2$,
 222 respectively. Thus the character a has rank 1, while t has rank 4. The
 223 characters c and g have rank 2 and 3, respectively.

224 It is important to notice that the description of the CDS approach pre-
 225 sented in this paper is slightly simplified compared to that introduced in [12].
 226 Specifically, in the original approach described in [12], use was made of the
 227 k -bounded position function, $\delta_k : \{1, \dots, n_c\} \rightarrow \{0, \dots, k - 1\}$, where k is a
 228 given threshold constant and $\delta_k(i)$ is the position (modulus k), i.e. $\delta_k(i) =$
 229 $[\delta(i) \bmod k]$, for each $i = 1, \dots, n_c$. Then the k -bounded-position sampled
 230 version of y , indicated by \dot{y} , is a numeric sequence, of length n_c defined as
 231 $\dot{y} = \langle \delta_k(1), \delta_k(2), \dots, \delta_k(n_c) \rangle$. Plainly we have $0 \leq \dot{y}[i] < k$, for each $1 \leq i \leq n_c$.

232 Although this allows to store each element of the sampled version of the
 233 text using only $\log(k)$ bits, in order to be able to retrieve the original i -th
 234 position $\delta(i)$, of the pivot character, from the i -th element of the k -bounded
 235 position sampled text \dot{y} , we need to maintain a *block-mapping table* τ which
 236 stores the indexes of the last positions of the pivot character in each k -block
 237 of the original text. Then, if we assume that the text y is divided in $\lceil n/k \rceil$
 238 blocks of length k , with the last block containing $(n \bmod k)$ characters,
 239 then we have $\tau[i] = \max(\{j : \delta(j) \leq ik\} \cup \{0\})$, for $1 \leq i \leq \lceil n/k \rceil$. Thus it
 240 is trivial to prove that $\tau[i] = j$ if and only if $\delta(j) \leq (ik)$ and $\delta(j + 1) > (ik)$.
 241 In addition, since the values in the block mapping τ are stored in a non
 242 decreasing order, i.e. $\tau[i] \leq \tau[i + 1]$, $\forall 0 \leq i \leq \lceil n/k \rceil$, it follows that
 243 $\delta(j) = (\tau[b] - 1)k + \dot{y}[j]$, where $b = \min\{i : \tau[i] \geq j\}$.

244 In practical cases, the choice of $k = 256$ would allow the sampled text to
 245 be stored using n_c bytes, plus the additional space of $4n/k$ bytes to maintain
 246 the block-mapping table.

247 In this paper, we simplified the representation of the CDS approach by
 248 dropping the τ table and keeping only the sampled text by storing each
 249 position with 4 bytes. On the one hand, this makes the approach more
 250 efficient from a practical point of view by avoiding the computation of the
 251 $\delta(i)$ value at each step of the algorithm. On the other hand, this increases
 252 the space required for storing the partial index. However, this disadvantage
 253 is largely mitigated by the extension of the approach to condensed alphabets,
 254 presented in the next section.

255 3.1. Extension to Condensed Alphabets

Let y be an input string, of length n , over an alphabet Σ of size σ . Given
 a constant parameter q , with $1 \leq q < n$, we define the condensed alphabet
 $\Sigma_y^{(q)}$, related to y , as

$$\{c \in \Sigma^q \mid c = y[i..i + q - 1] \text{ for some } 1 \leq i \leq n - q + 1\}$$

256 Roughly speaking $\Sigma_y^{(q)}$ is the set of all different substrings of length q (or
 257 q -grams) appearing in y . We define the q -condensed version of y as follows.

Definition 4 (q -Condensed Sequence). *Let y be a text of length n over an alphabet Σ of size σ and let $\Sigma_y^{(q)}$ be the condensed alphabet, related to y , for a given constant parameter q . We define the q -condensed version of the sequence y as the sequence, of length $n - q + 1$, of all consecutive (and overlapping) substrings of length q appearing in y . More formally*

$$y^{(q)} = \langle y[1..q], y[2..q+1], y[3..q+2], \dots, y[n-q+1..n] \rangle$$

Example 5. *Assume $y = \text{“agtagcgcagt”}$ is a DNA sequence of length 11, over the alphabet $\Sigma = \{a, c, g, t\}$. Then we have*

$$\begin{aligned} y^{(2)} &= \langle ag, gt, ta, ag, gc, cg, gc, ca, ag, gt \rangle \\ y^{(3)} &= \langle agt, gta, tag, agc, gcg, cgc, gca, cag, agt \rangle \\ y^{(4)} &= \langle agta, gtag, tagc, agcg, gcgc, cgca, gcag, cagt \rangle \end{aligned}$$

258 **Definition 5** (q -Characters Distance Sampling). *Let $C \subseteq \Sigma_y^{(q)}$ be the set of
 259 pivot characters, let $n_C \leq n$ be the number of occurrences of any pivot char-
 260 acter in the text $y^{(q)}$ and let δ be the position function of $y^{(q)}$. We define the
 261 q -characters distance function $\Delta^{(q)}$ as the distance between two consecutive
 262 occurrences of any pivot character in $y^{(q)}$, where $\Delta^{(q)}(i)$, for $1 \leq i \leq n_C - 1$, is
 263 the distance between the $(i+1)$ -th and the i -th occurrence of any occurrences
 264 of any pivot character in $y^{(q)}$.*

265 *Then the q -characters-distance sampled version of the sequence y is a
 266 numeric sequence of length $n_C - 1$, indicated by $\bar{y}^{(q)}$ and defined as*

$$\bar{y}^{(q)} = \langle \Delta^{(q)}(1), \Delta^{(q)}(2), \dots, \Delta^{(q)}(n_C - 1) \rangle. \quad (3)$$

Example 6. *As in the previous Example 5 assume $y = \text{“agtagcgcagtagta”}$ is a DNA sequence of length 15, over the alphabet $\Sigma = \{a, c, g, t\}$. If we suppose $q = 2$ and $C = \{\text{“ag”}\}$ is the set of pivot characters, then we have*

$$\begin{aligned} \dot{y}^{(2)} &= \langle 1, 4, 9, 12 \rangle \\ \bar{y}^{(2)} &= \langle 3, 5, 3 \rangle. \end{aligned}$$

Similarly, if we suppose $q = 3$ and $C = \{\text{“agt”}\}$ is the set of pivot characters, then we have

$$\begin{aligned} \dot{y}^{(3)} &= \langle 1, 9, 12 \rangle \\ \bar{y}^{(3)} &= \langle 8, 3 \rangle. \end{aligned}$$

267 In this paper we do not go into the way for a correct selection of the set
268 of pivot characters, and even we leave the details of an analysis about what
269 is the best subset to be chosen. However in our experimental evaluation (see
270 Section 6) we will show how it is enough to put a single character in the
271 set of pivot characters. Such a character is selected on the basis of its *rank*
272 *value*, where we remember that the rank of a character c corresponds to its
273 position in the alphabet Σ when we assume that all characters are sorted by
274 their frequencies inside the text (see Definition 3).

275 As we will note later, even the choice of the most frequent character of
276 the alphabet (be it a single character or a q -gram) is enough to obtain a good
277 efficacy, both in terms of search time and in terms of space used for storing
278 the partial index. In the following two sections we will show how to apply
279 sampling techniques in the two main scenarios, that of online searching and
280 that of offline searching.

281 4. Online Sampled String Matching

282 In this section we show how the sampled text-based approach can be
283 adopted to solve the online string matching problem. Specifically, in a first
284 phase, we briefly present how the OTS approach is used in this scenario.
285 Next we present the solution for online string matching based on the CDS
286 approach, also based on a condensed alphabet. Such algorithms make use
287 of an auxiliary string matching algorithm, used for searching the sampled
288 pattern on the sampled text, and they work well with most of the known
289 string matching algorithms. However, since the sampled patterns tend to be
290 short, we assume that the search phase is implemented using the Horspool
291 algorithm, which has been found to be fast in such setting. Such assumption
292 is the same as that adopted in the paper by Claude *et al.* [5]

293 4.1. Online Searching Using the OTS Approach

294 Claude *et al.* [5] presented a very efficient algorithm for online string
295 matching based on their OTS approach. Specifically, let y be the input text,
296 of length n , and let x be the input pattern, of length m , both over an alphabet
297 Σ of size σ . In addition let $\hat{\Sigma} \subset \Sigma$ be the sampled alphabet, and let \hat{y} be
298 the sampled text of length \hat{n} , containing all (and only) the characters of the
299 sampled alphabet.

The OTS algorithm constructs a sampled version of the input pattern,
 \hat{x} , of length \hat{m} during the searching phase. Such pattern is then searched in

the sampled text. Since \hat{y} contains partial information, for each candidate position i returned by the search procedure on the sampled text, the algorithm has to verify the corresponding occurrence of x in the original text. For this reason it uses information maintained in the table ρ to map positions of the sampled text to their corresponding positions in the original text. The position mapping ρ has size $\lfloor \hat{n}/h \rfloor$, where h is the *interval factor*, and is such that $\rho[i] = j$ if character $y[j]$ corresponds to character $\hat{y}[h \times i]$. More formally we have, for $1 \leq i \leq \lfloor \hat{n}/h \rfloor$

$$\rho[i] = j, \text{ if } y[j] \text{ is the } (h \times i)\text{-th occurrence in } y \text{ of any character of } \hat{\Sigma}$$

300 The value of $\rho[0]$ is set to 0. In their paper, on the basis of an accurate
301 experimentation, the authors suggest to use values of h in the set $\{8, 16, 32\}$.

302 Then, if the candidate occurrence position j is stored in the mapping
303 table, i.e if $\rho[i] = j$ for some $1 \leq i \leq \lfloor \hat{n}/h \rfloor$, the algorithm directly checks
304 the corresponding position in y for the whole occurrence of x . Otherwise, if
305 the sampled pattern is found in a position r of \hat{y} , which is not mapped in
306 ρ , the algorithm has to check the substring of the original text which goes
307 from position $\rho[r/h] + (r \bmod h) - \alpha + 1$ to position $\rho[r/h + 1] - (h - (r$
308 $\bmod h)) - \alpha + 1$, where α is the first position in x such that $x[\alpha] \in \hat{\Sigma}$.

309 Notice that, if the input pattern does not contain characters of the sam-
310 pled alphabet, i.e. $\bar{m} = 0$, the algorithm merely reduces to search for x in
311 the original text y .

312 **Example 7.** Suppose $y = \text{“abaacabdaacabcc”}$ is a text of length 15 over the
313 alphabet $\Sigma = \{a, b, c, d\}$. Let $\hat{\Sigma} = \{b, c, d\}$ be the sampled alphabet, by omitting
314 character “a”. Thus the sampled text is $\hat{y} = \text{“bcbdcbcc”}$. If we map every
315 $h = 2$ positions in the sampled text, the position mapping ρ is $\langle 5, 8, 12, 14 \rangle$.
316 To search for the pattern $x = \text{“acab”}$ the algorithm constructs the sampled
317 pattern $\hat{x} = \text{“cb”}$ and search for it in the sampled text, finding two occurrences
318 at position 2 and 5, respectively. We note that $\hat{y}[2]$ is mapped and thus it
319 suffices to verify for an occurrence starting at position 4, finding a match.
320 However position $\hat{y}[5]$ is not mapped, thus we have to search in the substring
321 $y[\rho(2) + 3 - 1.. \rho(3)]$, finding a match at position 10.

322 The real challenge in their algorithm is how to choose the best alphabet
323 subset to sample. Based on some analytical results, supported by an experi-
324 mental evaluation, they showed that it suffices in practice to sample the least

325 frequent characters up to some limit.⁵ are removed from the original alpha-
 326 bet. Under this assumption their algorithm has an extra space requirement
 327 which is only 14% of text size and is up to 5 times faster than standard online
 328 string matching on English texts.

329 4.2. Online Searching Using the CDS Approach

330 Let y be an input text of length n over an alphabet Σ of size σ , let $q > 1$
 331 and let $\Sigma_y^{(q)}$ be the condensed alphabet over Σ . In addition let $C \subseteq \Sigma^{(q)}$ be
 332 the set of pivot characters.

333 During the preprocessing phase the algorithm performs a scanning of the
 334 text y and builds the corresponding position sampled text $\dot{y}^{(q)}$.

335 Let now x be an input pattern of length m and let m_C be the number of
 336 occurrences of any pivot character in $x^{(q)}$. The searching phase can be then
 337 divided in three different subroutines, depending on the value of m_C . All
 338 searching procedures work using a filtering approach. The idea behind such
 339 searching procedures is to take advantage of the sampled text $\dot{y}^{(q)}$ computed
 340 during the preprocessing phase in order to quickly locate any candidate sub-
 341 string s of the original text which may include an occurrence of the pattern.

342 If such candidate substring s has length m then the algorithm simply
 343 performs a character-by-character comparison between the pattern and the
 344 substring. Otherwise if the candidate substring s has length greater than m ,
 345 then a searching procedure is called, based on a standard exact online string
 346 matching algorithm, for searching the pattern x in s .

347 In what follows we describe in details the three different searching proce-
 348 dures which are applied when $m_C = 0$, $m_C = 1$ and $m_C > 1$, respectively.

349 **Case 1:** $m_C = 0$

350 If the pattern contains no occurrence of any pivot characters, we have that
 351 m_C is equal to 0. Under this assumption the algorithm searches for the
 352 pattern x in all substrings of the original text which do not contain the pivot
 353 characters. Specifically such substrings are identified in the original text by
 354 the intervals $[\delta^{(q)}(i) + 1.. \delta^{(q)}(i + 1) + q - 2]$, for each $0 \leq i \leq n_C$, assuming
 355 $\delta^{(q)}(0) = 0$ and $\delta^{(q)}(n_C + 1) = n - q + 2$.

⁵According to their theoretical evaluation and their experimental results it turns out that, when searching on an English text, the best performance are obtained when the 13 most frequent characters

```

SEARCH-0( $x, \dot{y}^{(q)}, y, q$ )
1.  $m \leftarrow \text{len}(x)$ 
2.  $n_C \leftarrow \text{len}(\dot{y})$ 
3.  $\dot{y}^{(q)}[0] \leftarrow 0$ 
3.  $\dot{y}^{(q)}[n_C + 1] \leftarrow n - q + 2$ 
4. for  $i \leftarrow 1$  to  $n_C + 1$  do
5.     if  $(\dot{y}^{(q)}[i] - \dot{y}^{(q)}[i - 1] + q - 2 \geq m)$  then
6.          $l \leftarrow \dot{y}^{(q)}[i - 1] + 1$ 
7.          $r \leftarrow \dot{y}^{(q)}[i] + q - 2$ 
8.         search for  $x$  in  $y[l..r]$ 

```

Figure 1: The pseudocode of procedure SEARCH-0 for the sampled string matching problem, when no pivot character occurs in the input pattern x .

356 Specifically, for each $1 \leq i \leq n_C + 1$, the algorithm checks if the value
357 $\dot{y}^{(q)}(i) - \dot{y}^{(q)}(i - 1) + q - 2$ is greater or equal to m . In such a case the algorithm
358 searches for x in the substring of the text $y[\dot{y}^{(q)}[i - 1] + 1.. \dot{y}^{(q)}[i] + q - 2]$
359 using any standard string matching algorithm. Otherwise the substring is
360 skipped, since no occurrence of the pattern could be found at such position.
361 The pseudocode of procedure SEARCH-0 for the sampled string matching
362 problem, when no pivot character occurs in the input pattern x , is depicted
363 in Figure 1.

364 **Case 2:** $m_C = 1$

365 If the pattern x contains a single occurrence of any character of the set
366 C , then the length of the sampled version of the pattern is still equal to 0.
367 However also in this case the algorithm is able to efficiently take advantage of
368 the information precomputed in $\dot{y}^{(q)}$ using the positions of the pivot character
369 in $y^{(q)}$ as an anchor to locate all candidate occurrences of x .

370 Specifically, let α be the unique position in x which contains the pivot
371 character, i.e. we assume that $x[\alpha.. \alpha + q - 1] = c$ and that both $x[1.. \alpha - 1]$ and
372 $x[\alpha + 1.. m]$ do not contain any pivot character. Then, for each $0 \leq i \leq n_C - 1$,
373 the algorithm checks if the value $\dot{y}^{(q)}(i - 1) - \dot{y}^{(q)}(i - 2)$ is greater than $\alpha - 1$
374 and if the value $\dot{y}^{(q)}(i) - \dot{y}^{(q)}(i - 1)$ is greater than $m - \alpha$. In such a case
375 the algorithm merely checks if the substring of the text $y[\dot{y}^{(q)}[i - 1] - \alpha +$
376 $1.. \dot{y}^{(q)}[i - 1] - \alpha + m]$ is equal to the pattern. Otherwise the substring is
377 skipped. As before we assume that $\dot{y}(0) = 0$ and $\dot{y}(n_C + 1) = n + 1$. The last

```

SEARCH-1( $x, \dot{y}^{(q)}, y, q$ )
1.  $m \leftarrow \text{len}(x)$ 
2.  $n_C \leftarrow \text{len}(\dot{y})$ 
3.  $\alpha \leftarrow \min\{i : x^{(q)}[i] \in C\}$ 
4.  $\dot{y}^{(q)}[0] \leftarrow 0$ 
5.  $\dot{y}^{(q)}[n_C + 1] \leftarrow n - q + 2$ 
6. for  $i \leftarrow 1$  to  $n_C + 1$  do
7.     if ( $\dot{y}^{(q)}[i - 1] - \dot{y}^{(q)}[i - 2] > \alpha - 1$  and
            $\dot{y}^{(q)}[i] - \dot{y}^{(q)}[i - 1] > m - \alpha$ ) then
8.          $l \leftarrow \dot{y}^{(q)}[i - 1] - \alpha + 1$ 
9.          $r \leftarrow \dot{y}^{(q)}[i - 1] - \alpha + m$ 
10.        compare  $x$  and  $y[l..r]$ 

```

Figure 2: The pseudocode of procedure SEARCH-1 for the sampled string matching problem, when the pattern x contains a single occurrence of the pivot character.

378 alignment of the pattern in the text is verified separately at the end of the
379 main cycle. The pseudocode of procedure SEARCH-1 for the sampled string
380 matching problem, when the pattern x contains a single occurrence of the
381 pivot character, is depicted in Figure 2.

382 **Case 3:** $m_C \geq 2$

383 If the number of occurrences of any pivot character in C is greater than
384 1 then the algorithm uses the sampled text $\dot{y}^{(q)}$ to compute on the fly the
385 sampled version $\bar{y}^{(q)}$ of $y^{(q)}$ and use it to search for any occurrence of $\bar{x}^{(q)}$.
386 This is used as a filtering phase for locating in y any candidate occurrence.

387 First the character distance sampled version \bar{x} of x is computed. Then
388 the algorithm searches for \bar{x} in \bar{y} using any exact online string matching
389 algorithm. Notice that \bar{y} can be efficiently retrieved online from the sampled
390 text \dot{y} , using relation given in (2).

391 For each candidate occurrence i of \bar{x} located in \bar{y} , an additional procedure
392 must be run to check if such occurrence corresponds to a match of the whole
393 pattern x in y . For this purpose the algorithm checks if the substring of the
394 text $y[\dot{y}^{(q)}[i] - \dot{y}^{(q)}[0] .. \dot{y}^{(q)}[i] + m - 1]$ is equal to x , where $\dot{y}^{(q)}[0]$ is the position of
395 the first occurrence of the pivot character into the pattern. The pseudocode
396 of procedure SEARCH-2⁺ for the sampled string matching problem, when the
397 pattern x contains at least 2 occurrences of the pivot character, is depicted

```

SEARCH-2+( $x, \dot{y}^{(q)}, y, q$ )
1.  $m \leftarrow \text{len}(x)$ 
2.  $(\bar{x}^{(q)}, \bar{m}) \leftarrow \text{COMPUTE-DISTANCE-SAMPLING}(x, m, C)$ 
3. search for  $\bar{x}^{(q)}$  in  $\bar{y}^{(q)}$  :
4.   for each  $i$  such that  $\bar{x}^{(q)} = \bar{y}^{(q)}[i..i + \bar{m} - 1]$  do
5.      $l \leftarrow \dot{y}^{(q)}[i] - \dot{y}^{(q)}[0]$ 
6.      $r \leftarrow \dot{y}^{(q)}[i] + m - 1$ 
7.     compare  $x$  and  $y[l..r]$ 

```

Figure 3: The pseudocode of procedure SEARCH-2⁺ for the sampled string matching problem, when the pattern x contains at least 2 occurrences of the pivot character.

398 in Figure 3.

399 4.3. Complexity Issues

400 In this section we prove that, assuming an underlying auxiliary string
401 matching algorithm with a linear worst case and a $O(n \log(m)/m)$ average
402 case time complexity, the resulting algorithm based on Character Distance
403 Sampling over condensed alphabets achieves an optimal $O(n)$ time complex-
404 ity in the worst-case and a $O(n \log(m)/m)$ time complexity in the average
405 case. The following lemmas prove that procedures SEARCH-0, SEARCH-1
406 and SEARCH-2⁺, respectively, achieve, under suitable conditions, optimal
407 time complexity in both worst and average cases.

408 **Lemma 1.** *Let x and y be two strings of size m and n , respectively, over an*
409 *alphabet Σ of size $\sigma > 1$. Let $C \subseteq \Sigma^{(q)}$ be set of pivot characters and let $\dot{y}^{(q)}$*
410 *be position sampled version of the text y . Under the assumption of equiprob-*
411 *ability and independence of characters in Σ , the worst-case and average time*
412 *complexity of SEARCH-0 are $O(n)$ and $O(n \log_{\sigma} m/m)$, respectively.*

Proof. In our argumentation we refer to the pseudo-code reported in Fig.1. In order to evaluate the worst-case time complexity of SEARCH-0, we can notice that each substring of the text is scanned at least once in line 8, with no overlap. Thus if we use a linear algorithm to perform the standard search then it is trivial to prove that the whole searching procedure requires

$$T_{\text{wst}}^0(n) = O(m) + \sum_{i=1}^{n_c-1} O(\Delta^{(q)}(i)) = O(n).$$

413 where, the $O(n)$ term is related to the pre-processing of the pattern, while
 414 each $O(\Delta^{(q)}(i))$ term corresponds to the time required to process the sub-
 415 string $y^q[\delta(i)..\delta(i+1)]$.

Assuming that the underlying algorithm has an $O(n \log m/m)$ average time complexity, on a text of length n and a pattern of length m , we can express the expected average time complexity as

$$T_{\text{avg}}^0(n) = \sum_{i=1}^{n_c-1} O\left(\frac{\Delta^{(q)}(i) \log_{\sigma} m}{m}\right) = O\left(\frac{n \log_{\sigma} m}{m}\right).$$

416

□

417 **Lemma 2.** *Let x and y be two strings of size m and n , respectively, over*
 418 *an alphabet Σ of size $\sigma > 1$. Let $C \subseteq \Sigma^{(q)}$ be the set of pivot characters*
 419 *and let $y^{(q)}$ be the position sampled version of the text y . Under the as-*
 420 *sumption of equiprobability and independence of characters in Σ , the worst-*
 421 *case and average time complexity of the SEARCH-1 algorithm are $O(n)$ and*
 422 *$O(n \log_{\sigma} m/m)$, respectively.*

423 *Proof.* In our argumentation we refer to the pseudo-code reported in Fig.2. In
 424 order to evaluate the worst-case time complexity of the procedure, notice that
 425 each character could be involved in, at most, two consecutive checks in line
 426 10. Specifically any text position in the interval $[\delta(i-1)+1..\delta(i)-1]$ could be
 427 involved in the verification of the substrings $y[\delta(i-1)-\alpha+1..\delta(i-1)+m-\alpha]$
 428 and $y[\delta(i)-\alpha+1..\delta(i)+m-\alpha]$. Thus the overall worst case time complexity
 429 of the searching phase is $T_{\text{wst}}^1(n) = O(n)$.

In order to evaluate the average-case time complexity of the procedure, notice that the expected number of occurrences in y^q of the set of pivot characters is given by $\mathbf{E}(n_c) = n/\sigma$. Moreover, for any candidate occurrence of x in y , the number $\mathbf{E}(\text{insp})$ of expected character inspections performed by procedure VERIFY, when called on a pattern of length m , is given by

$$\mathbf{E}(\text{insp}) = 1 + \sum_{i=1}^{m-1} \left(\frac{1}{\sigma}\right)^i \leq \frac{\sigma}{\sigma-1}$$

Thus the average time complexity of the algorithm can be expressed by

$$\begin{aligned} T_{\text{avg}}^1(n) &= \mathbf{E}(n_c) \cdot \mathbf{E}(\text{insp}) = \\ &O\left(\frac{n}{\sigma}\right) \cdot O\left(\frac{\sigma}{\sigma-1}\right) = \\ &O\left(\frac{n}{\sigma-1}\right) \end{aligned}$$

430 obtaining the optimal average time complexity $O(n \log_\sigma m/m)$ for great
 431 enough alphabets of size $\sigma > (m/\log_\sigma m) + 1$, and for $k \geq \sigma$. \square

432 **Lemma 3.** *Let x and y be two strings of size m and n , respectively, over
 433 an alphabet Σ of size $\sigma > 1$. Let $C \subseteq \Sigma^{(a)}$ be the set of pivot characters
 434 and let $\dot{y}^{(q)}$ be the position sampled version of the text y . Under the as-
 435 sumption of equiprobability and independence of characters in Σ , the worst-
 436 case and average time complexity of the SEARCH-2⁺ algorithm are $O(n)$ and
 437 $O(n \log_\sigma m/m)$, respectively.*

438 *Proof.* In our argumentation we refer to the pseudo-code reported in Fig.3.
 439 In order to evaluate the worst-case time complexity of the algorithm in this
 440 last case notice that, if we use a linear algorithm to search \bar{y} for \bar{x} , the
 441 overall time complexity of the searching phase is $O(n_c + n_x m)$, where n_x
 442 is the number of occurrences of \bar{x} in \bar{y} . In the worst case it translates in
 443 $O(n_c m)$ worst case time complexity. However it is not difficult to suppose
 444 to implement procedure VERIFY based on a linear algorithm, as KMP, in
 445 order to remember all positions of the text which have been already verified,
 446 allowing the algorithm to run in overall $T_{\text{wst}}^{2+}(n) = O(n)$ worst-case time
 447 complexity.

In order to evaluate the average-case time complexity of the algorithm
 notice that time required for searching \bar{x} in \bar{y} is $O(n_c \log m_c/m_c)$. Moreover,
 observe that the number of verification is bounded by the expected number
 of occurrences of the pivot character in y , thus, following the same line of
 Theorem 2, the overall average time complexity of the verifications phase is
 $O(n/(\sigma - 1))$. Thus the average time complexity of the algorithm can be
 expressed by

$$T_{\text{avg}}^{2+}(n) = O\left(\frac{n_c \log m_c}{m_c}\right) + O\left(\frac{n}{\sigma - 1}\right)$$

448 obtaining the optimal average time complexity $O(n \log_\sigma m/m)$ for great enough
 449 values of σ , such that $\sigma \geq (m/\log_\sigma m) + 1$. \square

450 5. Offline Sampled String Matching

451 In this section we describe an approach to indexed searching which makes
 452 use of a suffix array constructed over the sampled version of the text. In our
 453 evaluation we have chosen to use the suffix array [30] as a reference point since
 454 it can be counted among the most efficient standard solutions to the offline

455 string matching problem and because this solution was previously adopted
456 by Claude *et al.* [5] for comparison with their approach to text sampling.

457 We report that the suffix array data structure can be improved in various
458 ways. For instance it can be effectively compressed with relative Lempel-Ziv
459 (RLZ) dictionary compression, in such a way that arbitrary sub-arrays can be
460 rapidly decompressed, thus facilitating compressed indexing [34, 35]. Among
461 other possibilities we also report that a Suffix Array can be enhanced by rep-
462 resenting a sequence of integers using Fibonacci encodings, thereby reducing
463 the space requirements while retaining the searching functionalities [1]. In
464 addition a Suffix Array can be improved in efficiency in various ways [25].
465 However, these improvements are beyond the scope of this work, which in-
466 tends to verify how a sampled text-based approach can improve the search
467 efficiency of offline approaches, although it is possible to imagine that differ-
468 ent data structures can achieve different improvements.

469 We remember that a suffix array is a sorted array of all suffixes of a
470 string. Such data structure has been introduced by Manber and Myers in
471 1990 [30] as a simple, space efficient alternative to suffix trees [9]. It has been
472 extensively studied in the last three decades and in 2016 Li, Li and Huo [27]
473 gave the first in-place $O(n)$ -time construction algorithm that is optimal both
474 in time and space, where in-place means that the algorithm only needs $O(1)$
475 additional space beyond the input string and the output suffix array.

476 Formally, given a text y of length n , the suffix array s_y of y is defined
477 to be an array of integers providing the starting positions of suffixes of y
478 in lexicographical order. This means that $s_y[i]$ contains the starting posi-
479 tions of the i -th smallest suffix in y and thus for all $1 \leq i \leq n$, we have
480 $y[s_y[i-1]..n] < y[s_y[i]..n]$.

481

482 **Example 8.** Let $y = \text{“agaacgcagtata”}$ be a text of length 13, over the alpha-
483 bet $\Sigma = \{a, c, g, t\}$. The suffix array, s_y , contains the starting positions of all
484 suffixes of y , arranged in lexicographical order. Specifically we have:

485

	$s_y[1]$	$= 12$	$\rightarrow \langle a \rangle$
	$s_y[2]$	$= 2$	$\rightarrow \langle aacgcagtata \rangle$
	$s_y[3]$	$= 3$	$\rightarrow \langle acgcagtata \rangle$
	$s_y[4]$	$= 0$	$\rightarrow \langle agaacgcagtata \rangle$
	$s_y[5]$	$= 7$	$\rightarrow \langle agtata \rangle$
	$s_y[6]$	$= 10$	$\rightarrow \langle ata \rangle$
486	$s_y[7]$	$= 6$	$\rightarrow \langle cagtata \rangle$
	$s_y[8]$	$= 4$	$\rightarrow \langle cgcagtata \rangle$
	$s_y[9]$	$= 1$	$\rightarrow \langle gaacgcagtata \rangle$
	$s_y[10]$	$= 5$	$\rightarrow \langle gcagtata \rangle$
	$s_y[11]$	$= 8$	$\rightarrow \langle gtata \rangle$
	$s_y[12]$	$= 11$	$\rightarrow \langle ta \rangle$
	$s_y[13]$	$= 9$	$\rightarrow \langle tata \rangle$

487 The time complexity needed to build suffix array is $O(n^2 \log(n))$ if an
488 $O(n \log(n))$ algorithm is used for sorting the array of all suffixes. However,
489 there are many efficient algorithms to build suffix array [22]. Once the suffix
490 array is built, it is possible to search a pattern using the suffix array by a
491 binary search in $O(n \log(n))$ time. However it has been proved that we can
492 report all *occ* occurrences of a pattern in a text in $O(m + \log(n) + occ)$ [2].
493 In the following two subsections we show how the suffix array based solution
494 can be adapted to the OTS and CDS approaches presented in this paper.

495 5.1. Offline Searching Using the OTS Approach

496 To turn the sampling approach into an index, Claude *et al.* use a suffix
497 array to index the sampled positions of the text. When constructing the suffix
498 array, only suffixes starting with a sampled character will be considered, but
499 the sorting will still be done considering the full suffixes. The resulting
500 sampled suffix array is like the suffix array of the original text where suffixes
501 starting with unsampled characters have been omitted.

502 **Example 9.** Let again $y = \text{“agaacgcagtata”}$ be a text of length 13, over the
503 alphabet $\Sigma = \{a, c, g, t\}$. The sampled suffix array, $s_{\hat{g}}$, contains the starting
504 positions of all suffixes of y beginning with a character with the sampled alpha-
505 bet, arranged in lexicographical order. Specifically, assuming that $\hat{\Sigma} = \{a, g\}$,
506 we have:

507

$s_{\hat{y}}[1] = 12 \rightarrow \langle a \rangle$
 $s_{\hat{y}}[2] = 2 \rightarrow \langle aacgcagtata \rangle$
 $s_{\hat{y}}[3] = 3 \rightarrow \langle acgcagtata \rangle$
 $s_{\hat{y}}[4] = 0 \rightarrow \langle agaacgcagtata \rangle$
508 $s_{\hat{y}}[5] = 7 \rightarrow \langle agtata \rangle$
 $s_{\hat{y}}[6] = 10 \rightarrow \langle ata \rangle$
 $s_{\hat{y}}[7] = 1 \rightarrow \langle gaacgcagtata \rangle$
 $s_{\hat{y}}[8] = 5 \rightarrow \langle gcagtata \rangle$
 $s_{\hat{y}}[9] = 8 \rightarrow \langle gtata \rangle$

509 Search on the sampled suffix array is carried out as follows. Given a
510 pattern x the algorithm finds the first sampled character of the pattern.
511 Assume such character is at index j of x . The pattern is then partitioned into
512 the prefix $x[1..j-1]$ and the suffix starting with the first sampled character
513 $x[j..m]$. The algorithm then searches the sampled suffix array for the suffix
514 of the pattern like in an ordinary suffix array. Each candidate occurrence
515 located by this search will then be verified by comparing the prefix $x[1..j-1]$
516 against the text. Observe that the OTS suffix array can be used for searching
517 a text only for patterns that contain at least one sampled character.

518 5.2. Offline Searching Using the CDS Approach

519 The algorithm we propose is divided into two phases: a first *preprocessing*
520 *phase* which consists in the construction of a sampled version, $s_{\hat{y}}$, of the suffix
521 array and a *searching phase* which is used to search any pattern x of length
522 m in y making use of the suffix array $s_{\hat{y}}$ and the sampled text \hat{y} . We notice
523 that, as it happens in any offline string matching solution, the preprocessing
524 phase is performed only once for the construction of the partial index, while
525 the searching phase can be run for an indeterminate number of queries. We
526 notice also that the algorithm must maintain the original text y , the sampled
527 version of the text \hat{y} and the corresponding suffix array $s_{\hat{y}}$.

528 We are now ready to describe the preprocessing and the searching phase
529 of our new proposed algorithm.

530

531 As before, let y be an input text of length n over an alphabet Σ of size
532 σ and let $C \subseteq \Sigma$ be the set of pivot characters. During the preprocessing
533 phase the algorithm builds and stores the position sampled text \hat{y} of y . This
534 requires $O(n)$ -time and $O(n_c)$ -space, where n_c is the number of occurrences

535 of any pivot character in y . Subsequently a suffix array of \bar{y} is constructed
 536 on the fly using information maintained in \dot{y} .

537 As a consequence, when constructing the suffix array of \bar{y} , the algorithm
 538 takes into account only suffixes beginning with a pivot character in the origi-
 539 nal text, drastically reducing the space requirement for maintaining the whole
 540 index.

541 **Definition 6** (CDS Suffix Array). *Let y be a text of length n , let $C \subseteq \Sigma$
 542 be the set of pivot characters and let n_C be the number of occurrences of
 543 any $c \in C$ in the input text y . Let $\delta : \{1, \dots, n_C\} \rightarrow \{1, \dots, n\}$ be the position
 544 function and let \dot{y} be the position sampled version of y .*

545 *The CDS suffix array $s_{\bar{y}}$ of y is defined to be an array of all index posi-
 546 tions i , with $1 \leq i \leq n_C - 1$ such that $\bar{y}[s_{\bar{y}}[i] - 1..n_C - 1] < \bar{y}[s_{\bar{y}}[i]..n_C - 1]$.*

547

548 **Example 10.** *Let $y = \text{“agaacgcagtata”}$ be a text of length 13, over the al-
 549 phabet $\Sigma = \{a, c, g, t\}$. Let $C = \{a\}$ be the set of pivot characters. Thus the
 550 position sampled version of y is $\dot{y} = \langle 1, 3, 4, 8, 11, 13 \rangle$, while the character
 551 distance sampled version of y is $\bar{y} = \langle 2, 1, 4, 3, 2 \rangle$.*

552 *The sampled suffix array, $s_{\bar{y}}$, contains all positions in y , starting with a
 553 pivot character $c \in C$, arranged in lexicographical order with respect to the
 554 suffixes of \bar{y} . Specifically we have:*

555

$$\begin{aligned}
 s_{\bar{y}}[1] &= 1 && \rightarrow \langle 1, 4, 3, 2 \rangle \\
 s_{\bar{y}}[2] &= 4 && \rightarrow \langle 2 \rangle \\
 556 \quad s_{\bar{y}}[3] &= 0 && \rightarrow \langle 2, 1, 4, 3, 2 \rangle \\
 s_{\bar{y}}[4] &= 3 && \rightarrow \langle 3, 2 \rangle \\
 s_{\bar{y}}[5] &= 2 && \rightarrow \langle 4, 3, 2 \rangle
 \end{aligned}$$

557

558 *Thus, $s_{\bar{y}}[0] = 1$ indicates that the smallest suffix, in lexicographical order
 559 relative to \bar{y} , begins at position $\dot{y}[1]$ in y (it is $\bar{y}[1..4] = \langle 1, 4, 3, 2 \rangle$). Similarly
 560 $s_{\bar{y}}[1] = 4$ indicates that the second suffix begins at position $\dot{y}[4]$ in y .*

561 During the searching phase the algorithm uses the suffix array of the
 562 sampled text $s_{\bar{y}}$ as an index to quickly locate every occurrences of a sampled
 563 pattern \bar{x} in \bar{y} . Each of these occurrences is treated as a candidate occurrence
 564 of x in y , and as such it will be verified by a comparison procedure.

565 The searching algorithm works as a standard search on a suffix array. It is
 566 based of the fact that finding every occurrence of the pattern \bar{x} is equivalent

567 to find every suffix in \bar{y} that begins with the \bar{x} . Thanks to the lexicographical
 568 ordering of the suffix array, all such suffixes are grouped together and can
 569 be found efficiently with a single binary search, which locates the starting
 570 position of the interval.

571 For the sake of completeness we observe that both the OTS and the
 572 CDS suffix arrays resemble a sparse suffix array [21], which indexes regularly
 573 sampled text positions. However, such data structure only need to make one
 574 search of the sampled pattern, while using a sparse suffix array h searches
 575 are needed if the suffix array indexes every h -th position. The drawback of
 576 such data structures is that they can only be used for patterns that contain
 577 at least one sampled character, whereas the sparse suffix array can be used if
 578 the pattern length is at least q . The variance of the search time when using
 579 the sampled suffix array is also larger than when using a sparse suffix array
 580 because in the sampled suffix array we have much less control over the length
 581 of the string that is used in the suffix array search.

582 5.3. Complexity Issues

583 In order to compute the time complexity needed for searching a pattern
 584 x , of length m , in a text y , of length n , we assume that n_C is the number
 585 pivot characters appearing in y . Then finding the first position of a sampled
 586 pattern \bar{x} of length m_C in a suffix array $s_{\bar{y}}$ of length n_C takes $O(m_C \log n_C)$ -
 587 time [2] while finding the set of all ρ occurrences of \bar{x} in \bar{y} takes $O(\rho)$ -time.
 588 Since each occurrence must be verified in the original text we need $O(m\rho)$
 589 additional time for the verification phase. The overall time complexity of the
 590 searching algorithm is then $O(m_C \log(n_C) + m\rho)$.

591 It is important to notice that such complexity leads to a worst-case sce-
 592 nario that cannot be compared with the $O(m \log(n))$ -time complexity ob-
 593 tained by suffix arrays. Indeed, the verification phase component $O(m\rho)$,
 594 which is not required in standard suffix arrays, may in many cases be domi-
 595 nant over the search phase component $O(m_C \log(n_C))$.

596 If we assume, for instance that $y = (abc)^n$ and $x = (abc)^{m-1}(acb)^m$ and
 597 the set of pivot characters is $C = \{a\}$, the sampled suffix array reports all
 598 text positions as candidate occurrences and the verification needs to be run
 599 n time. This scenario leads to an overall worst-case time complexity equal
 600 to $O(nm)$ overall.

601 6. Experimental Results

602 In this section we present experimental results obtained by comparing
603 the new proposed sampling approaches (using values of q ranging between 2
604 and 4) against the standard Character Distance Sampling (CDS) approach
605 (obtained with q set to 1). Experimental evaluations were conducted in both
606 the online and offline scenarios.

607 We compare our approaches against the Occurrence Text Sampling (OTS)
608 approach and against the standard solution for which no text sampling is
609 included.⁶ In order to conduct a comparison as fair as possible we also im-
610 plemented the OTS approach using q -grams, for values of q ranging between
611 2 and 4

612 All algorithms have been implemented using the C programming lan-
613 guage, and have been tested using a variant of the SMART⁷ tool [10] prop-
614 erly tuned for testing string matching algorithms based on a text-sampling.
615 Tests have been executed on a MacBook Pro with 4 Cores, a 2.7 GHz In-
616 tel Core i7 processor, 16 GB RAM 2133 MHz LPDDR3, 256 KB of L2
617 Cache and 8 MB of Cache L3. The code of the algorithms used to per-
618 form the experiments presented in this section is available online at [https:](https://www.dmi.unict.it/faro/SAMPLING/)
619 [//www.dmi.unict.it/faro/SAMPLING/](https://www.dmi.unict.it/faro/SAMPLING/).

620 All algorithms, for both the online and offline scenarios, have been tested
621 on two 100MB text buffers containing a real biological sequence and a natural
622 language text in the English language. Specifically the biological sequence
623 is a collection of newline-separated gene DNA sequences (without descrip-
624 tions, just the bare DNA code) obtained from files 01hg10 to 21hg10,
625 plus 0xhg10 and 0yhg10, from Gutenberg Project. Each of the 4 bases is
626 coded as an uppercase letter A, G, C, T, with few occurrences of other special
627 characters. The natural language text buffer is the concatenation of English
628 text files selected from `etext02` to `etext05` collections of Gutenberg Project,
629 where the headers related to the project have been deleted so as to leave just
630 the real text. Both sequences are available for download in the PIZZA&CHILI
631 Corpus (<http://pizzachili.dcc.uchile.cl>).

632 In the experimental evaluation (for both online and offline searching),

⁶The standard solutions taken as a reference point are the Boyer-Moore-Horspool al-
gorithm for the online scenario and the search on suffix arrays for the offline scenario.

⁷The SMART tool is available online for download at <http://www.dmi.unict.it/~faro/smart/> or at <https://github.com/smart-tool/smart>.

633 patterns of length m were randomly extracted from the sequences, with m
634 ranging over the set of values $\{2^i | 3 \leq i \leq 8\}$. For each value of m , the mean
635 over the running times (expressed in hundredths of seconds) of 1000 runs has
636 been reported.

637 In our implementations we selected the pivot character on the basis of its
638 *rank value*, where we remember that the rank of a character c is the position
639 of c in the alphabet Σ , if we assume that all characters are sorted by their
640 frequencies inside the text (see Definition 3).

641 Then we evaluated the behaviour of our algorithms for different values
642 of the rank r of the selected pivot character and specifically for r ranging
643 between 1 (the most frequent character) and 16. Observe that if σ is the size
644 of the original alphabet Σ , then σ^q is the size of the condensed alphabet $\Sigma^{(q)}$.
645 As a consequence, in the case of experimental tests on genome sequences and
646 $q = 1$, the value of the rank r is limited in the range between 1 and 4, since
647 4 is the size of the alphabet. We underline also that, in the case of the OTS
648 approach, the value of the rank r refers to variations of the size of the set of
649 sampled characters. Also in this case r ranges from 2 to 16.

650 6.1. Space Requirements

651 In the context of text-sampling string-matching space requirement is one
652 of the most significant parameter to take into account. It indicates how much
653 additional space, with regard to the size of the original input sequences, is
654 required by a given solution to solve the problem.

655 Text-sampling algorithms require to store the whole text together with
656 the additional sampled-text which is used to speed-up the searching phase.
657 Although sampled texts have the good property to allow a direct access to
658 the input text (when they are scanned sequentially), to be of any practical
659 interest they should require as little extra space as possible.

660 Fig. 4 and Fig. 5 show the space consumption of the newly proposed text-
661 sampling approaches, in the case of a genome sequence and an English text,
662 respectively. Data are reported for different values of q in terms of percentage
663 of memory used, in comparison with the original text size. We recall that,
664 in the case of CDS, memory space consumption is plotted on variations of
665 the rank of the pivot character. As expected, the function which describes
666 memory requirements shows a decreasing trend while the rank of the pivot
667 character increases. Similarly space consumption drastically decreases when
668 the size of q increases.

SPACE CONSUMPTION FOR A GENOME SEQUENCE

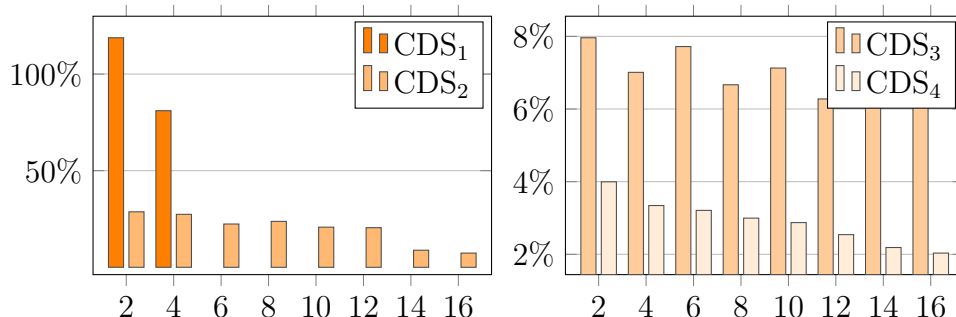


Figure 4: Space consumption of CDS approaches on a genome sequence, for different pivot characters with rank ranging from 2 to 16 and for different values of the parameter q , ranging from 1 to 4. Data are reported in terms of percentage of memory used relative to the original text size.

669 Data reported in Fig. 4, related to the sampling of a genome sequence,
 670 show that, when compared against the standard sampling approach (obtained
 671 with $q = 1$), the benefit in space consumption obtained by the approaches
 672 based on condensed alphabets is impressive. Specifically the gain for CDS
 673 ranges from 72% (for $r = 1$ and $q = 2$) to 95% (for $r = 16$ and $q = 4$).
 674 In addition we can observe a sensible gain in the space consumption also in
 675 comparison with the OTS algorithms implemented using condensed alphabets.
 676

677 Data reported in Fig. 5, related to the English text, show that, when
 678 compared against the standard sampling approach (obtained with $q = 1$), the
 679 benefit in space consumption obtained by the approaches based on condensed
 680 alphabets is even more advantageous. Specifically the gain for CDS ranges
 681 from 90% (for $r = 1$ and $q = 2$) to 98.8% (for $r = 16$ and $q = 4$). Also
 682 in this case, the CDS approach shows a significant reduction in the space
 683 consumption when compared with the OTS approach implemented using
 684 condensed alphabets.

685 For the sake of completeness we would like to point out that standard al-
 686 gorithms for the online string matching problem require an amount of space
 687 which is, in general, proportional to the length of the pattern and/or to the
 688 size of the alphabet. In this particular case (a 5MB text buffer) the Boyer-
 689 Moore-Horspool algorithm requires only 1.24 KB of memory for implement-
 690 ing the occurrence heuristic (equivalent to a $O(\sigma)$ -space complexity), while

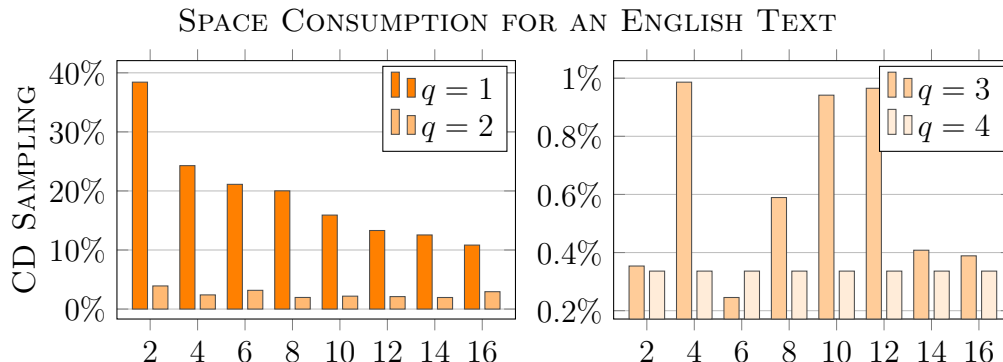


Figure 5: Space consumption of CDS approaches on an English text, for different pivot characters with rank ranging from 2 to 16 and for different values of the parameter q , ranging from 1 to 4. Data are reported in terms of percentage of memory used relative to the original text size.

691 some among the most effective algorithms (for instance $WFRq$ [4], $SKIPq$
 692 [8]) are implemented by means of a hash table of size 65536, requiring 0.2
 693 MB of additional space. Thus it turns out that, under particular condi-
 694 tions (texts of moderate lengths), the practical space requirements of our
 695 proposed sampling algorithms are comparable with those of standard online
 696 string matching solutions.

697 6.2. Online Searching

698 In this section we compare the different text-sampling approaches in terms
 699 of online searching times. In this context we refer to the *searching time* as
 700 the time needed to perform the searching of the pattern on both sampled
 701 and original texts, including any preprocessing of the underlying algorithm.
 702 However, in our analysis the searching time doesn't include the preprocessing
 703 time needed to construct the partial index.

704 Following the same lines of previous papers on sampled string matching
 705 [5, 12] we tested all sampling solutions in combination with the Boyer-Moore-
 706 Horspool (HOR) algorithm [20] for the implementation of the underlying
 707 standard searching procedure. As a consequence, in our comparison we also
 708 included the Boyer-Moore-Horspool string matching algorithm (in its stan-
 709 dard implementation) in order to understand how much the proposed sam-
 710 pling approach contributes to speed-up a standard online string matching

711 solution.⁸

In our experimental results we also included the best running time obtained by OTS solutions implemented with condensed alphabets. Specifically if $\text{OTS}_{(q,r)}$ is the searching speed of the OTS algorithm implemented using q -grams and rank r , we compute $\text{OTS}_{(q,r)}$ as

$$\text{OTS}_{(q,r)} = \max_{\substack{1 \leq q \leq 4 \\ 1 \leq r \leq 16}} (\text{OTS}_{(q,r)})$$

712 Observe that the original OTS approach is obtained by setting $q = 1$. Our
713 experiments showed that the values that yield the best results are $q = 3$ and
714 $r = 12$ for searching genomic sequences, and $q = 4$ and $r = 12$ for searching
715 English texts. Thus we have reported only these values in the graphs.

716 In addition, for the sake of completeness, we also included in our ex-
717 perimental results three among the most efficient algorithm recently intro-
718 duced for the exact online string matching problem. Specifically we included
719 the Weak-Factor-Recognition algorithm [4] (WRF), the Brackward-Range-
720 Automaton-Matcher [14] (BRAM) and the Skip-Search algorithm [8]. All
721 algorithms have been implemented in several variants using q -grams, for val-
722 ues of q ranging from 1 to 8. Here again, we have only reported the results
723 obtained with the best variant for each algorithm.

724 Fig. 6 and Fig. 7 show the resulting searching times of all tested algo-
725 rithms when they were used for searching on a genome sequence and on an
726 English text, respectively. Results are expressed in terms of searching speed,
727 reported in Gigabytes per second (GB/s).

728 In general, the search speeds achieved by the new variants are extremely
729 high and their advantage over the performance of standard algorithms is
730 impressive. This is due to the fact that the main loop of the search phase
731 iterates over the sampled text, the length of which, as we have seen above,
732 is much shorter than that of the original text. Added to this is the fact that
733 the number of candidate occurrences is very low, especially for long patterns.

734 From experimental results on a genome sequence (Fig. 6) it turns out that
735 in all cases the best results are obtained by the variants based on condensed

⁸Although there exists many other searching algorithms able to show better practical performances on biological data (see for instance [4, 8]) this kind of comparison goes beyond the objectives of this paper. We expect that the proposed approach is able to enhance the performances of different string matching algorithms with different, though similar, rates.

ONLINE SEARCHING ON GENOME SEQUENCE

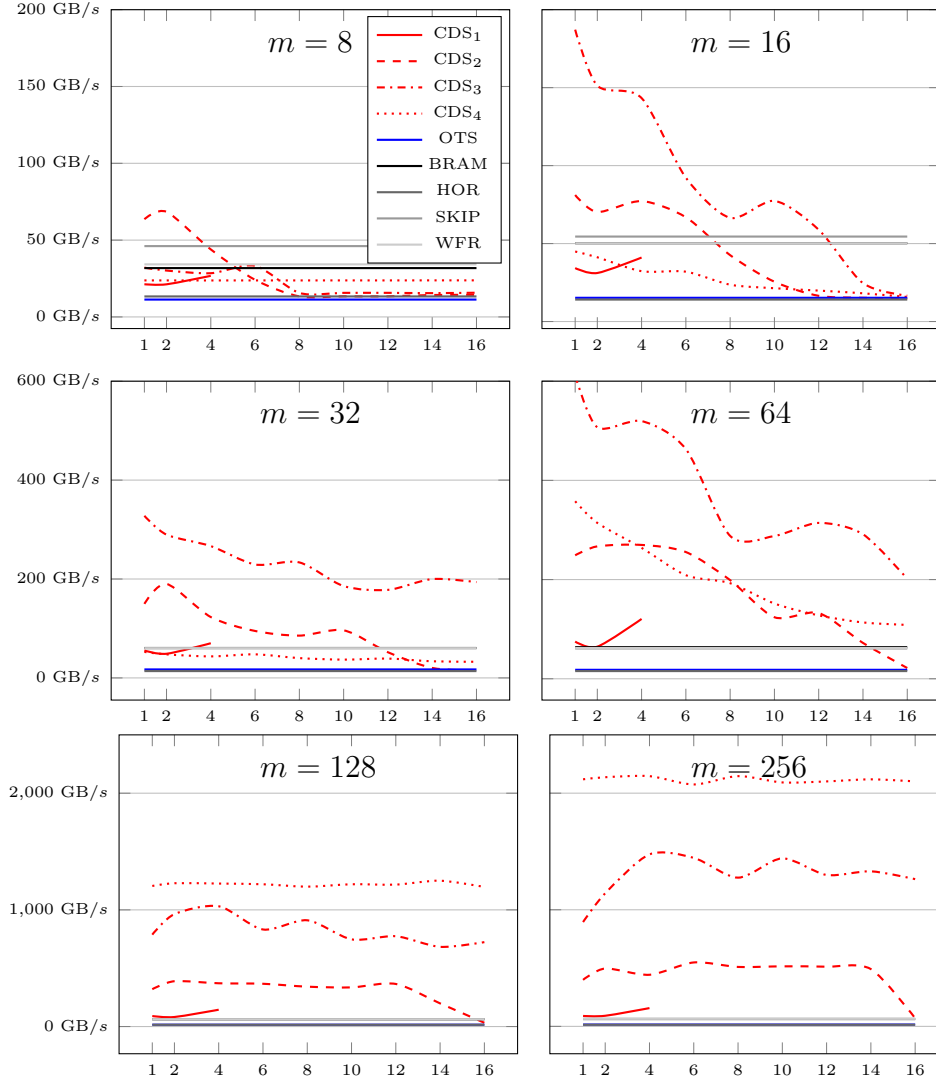


Figure 6: Searching speed on a genome sequence. Red lines represent the CDS_q algorithm implemented with $1 \leq q \leq 4$, the solid gray-tones lines represent the standard algorithms while the blue solid line represents the best searching time of the OTS solution implemented with q -grams. The x axis represents the rank r of the pivot character in the case of the sampling algorithms, with $1 \leq r \leq 16$.

ONLINE SEARCHING ON ENGLISH TEXT

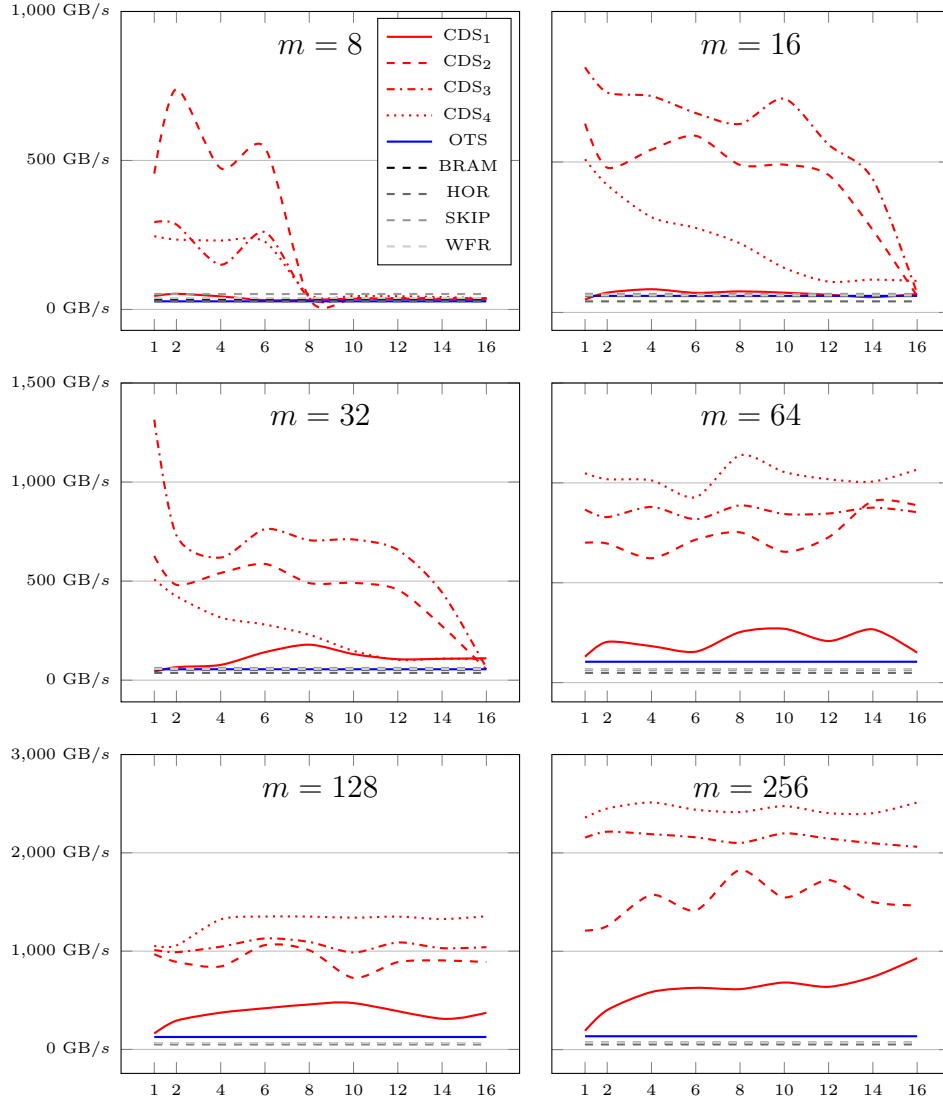


Figure 7: Searching speed on a natural language text. Red lines represent the CDS_q algorithm implemented with $1 \leq q \leq 4$, the solid gray-tones lines represent the standard algorithms while the blue solid line represents the best searching time of the OTS solution implemented with q -grams. The x axis represents the rank r of the pivot character in the case of the sampling algorithms, with $1 \leq r \leq 16$.

736 alphabets. Specifically, in the case of short patterns ($m = 8$) the best running
737 times are obtained by the variant with $q = 2$, but as the length of the pattern
738 increases, the variants with higher q values prove to be faster and faster until
739 the variant with $q = 4$ proves to be the fastest from $m \geq 128$.

740 As might be expected, search speeds also increase as the pattern length
741 increases, from speeds of just under 100 GB/s for $m = 8$ to speeds of just
742 over 2.000 GB/s for $m = 256$.

743 When using a value of q greater than 1, the speed up obtained by CDS
744 is always greater than 50% and reaches the value of 90% under suitable
745 conditions, i.e. for $q = 4$ and long patterns.

746 Observe that the behaviour of algorithms based in CDS follow a decreas-
747 ing trend for increasing rank values. Thus in most cases the better choice is
748 to use the most frequent element as the pivot character. Observe indeed that,
749 when the rank of the pivot character is greater than a given threshold, the
750 performances of the CDS algorithms based on q -grams sensibly degrades.
751 Specifically this threshold is approximately equal to 6 for short patterns
752 ($q \leq 3$ and $m = 8$), while it increases as the pattern gets longer or for greater
753 values of q .

754 Similarly, from experimental results on an English text (Fig. 7) it turns
755 out that the variants based on condensed alphabets obtain the best results
756 in all cases. Again, the best choice for short patterns is $q = 2$, but as the
757 length of the pattern increases, the variants with higher q values prove to be
758 faster and faster until the variant with $q = 4$ proves to be the fastest from
759 $m \geq 64$. The maximum speed reached by such solution is very close to
760 3.000 GB/s. When using a value of q greater than 1, the speed up obtained
761 by CDS is always greater than 50% and reaches the value of 98% under
762 suitable conditions, i.e. for $q = 4$ and long patterns.

763 We observe that on natural language texts and long patterns the be-
764 haviour of algorithms based in CDS follow a slightly increasing trend for
765 increasing rank values, but, in general, the search speed obtained from the
766 different values of the rank is comparable. As in the case of genome sequence,
767 in the case of short patterns, when the rank of the pivot character is greater
768 than a given threshold, the performances of the CDS algorithms based on
769 q -grams sensibly degrades.

770 Going into details of the improvement in terms of running times we ob-
771 serve that the original CDS approach ($q = 1$) leads to improvements which
772 are in percentage between 74% (in the case of short patterns) and 77% (in

773 the case of long patterns) if compared with the underlying standard string
774 matching algorithm. The new CDS algorithms based on condensed alpha-
775 blets give instead much more evident improvements which range from 96%
776 (for short patterns) and 99.6% (in the case of long patterns) compared with
777 the same algorithm. This improvements translate into a gain up to 70% for
778 short patterns and up to 96% in the case of long patterns, if compared with
779 OTS approach.

780 6.3. Offline Searching

781 In this section we compare the different text-sampling approaches in the
782 offline scenario. In this context we refer to the *searching time* as the time
783 needed to perform the searching of the pattern on the text-index. In our
784 analysis the searching time doesn't include the preprocessing time needed to
785 construct the index.

786 Following the same lines of [5] we tested all sampling solutions using a
787 modified Suffix Array, as described in Section 5. Thus, in our comparison
788 we also included the original Suffix Array algorithm (STD), in its standard
789 implementation, in order to understand how much the proposed sampling
790 approaches contribute to speed-up a standard offline searching solution.

791 We mention that Ferragina and Manzini [15] showed that it is possible
792 to search a pattern x of length m backwards in the suffix array of y without
793 storing it. A backward search means that we first search for the substring
794 $x[m..m]$, then for the substring $x[m-1..m]$, and so on, until the whole pattern
795 x is found. In the computer science literature, any data structure that allows
796 to search a pattern x backwards in the (conceptual) suffix array of a text y
797 is called an FM-index of y . In our experimental results, however, we have not
798 used the backward-search technique, limiting ourselves to a simple binary
799 search within the suffix array constructed on the text y .

800 Table. 1 and Table. 2 show the resulting performance of all tested algo-
801 rithms when they were used for searching on a genome sequence and on an
802 English text, respectively. Results are reported in terms of searching speed,
803 expressed in number of queries per second (QR/s).

804 For an easier reading of the results, we have listed in the tables, for
805 each algorithm, only the results of the variants with the best performance.
806 Specifically for OTS we reported the results with $q = 4$ and $r = 12$, in the
807 case of genomic sequences, and $q = 3$ and $r = 12$, in the case of natural
808 language texts. In the case of CDS we reported the results for pairs of values
809 $(q, r) \in \{(1, 2), (2, 8), (3, 10), (4, 8)\}$.

m	STD	OTS _(4,12)	CDS _(1,2)	CDS _(2,8)	CDS _(3,10)	CDS _(4,8)
8	152	161	188	191	201	198
16	154	172	183	197	213	227
32	161	183	201	217	216	234
64	163	199	218	234	241	268
128	168	202	233	247	246	299
256	181	216	257	281	293	322

Table 1: Offline searching on a genome sequence. Values are reported in thousands of queries per second. We used values q ranging from 1 to 4, values of r ranging from 1 to 16 and pattern lengths m ranging from 8 to 256. Best results have been bold-faced.

m	STD	OTS _(3,12)	CDS _(1,8)	CDS _(2,8)	CDS _(3,10)	CDS _(4,8)
8	198	212	244	244	247	261
16	212	224	257	254	271	267
32	219	222	271	278	284	299
64	218	238	321	338	337	355
128	219	243	322	343	356	385
256	224	278	366	368	382	401

Table 2: Offline searching on a natural language text. Values are reported in thousands of queries per second. We used values q ranging from 1 to 4, values of r ranging from 1 to 16 and pattern lengths m ranging from 8 to 256. Best results have been bold-faced.

810 From the experimental results it turns out that the standard solution
811 based on a suffix array offers performances of about 200K queries per second,
812 while the solutions based on the OTS approach oscillate between 200K and
813 300K queries per second, proposing a search speed mildly faster than the
814 previous one. The solutions based on the CDS approach, on the other hand,
815 offer significantly better performance, oscillating between 250K queries per
816 second ($q = 1$ and short patterns) and 400K queries per second ($q = 4$ and
817 for almost all values of m). These last solutions offer performances that are
818 therefore between 1.5 and 1.8 times faster than the standard solution.

819 Finally, we note that, as might be expected, there is no significant vari-
820 ation in the results for different values of m . In fact, although there is a
821 difference between the results, going from patterns of length 8 to patterns
822 of length 256, the gain obtained is at most just over 50%. This is due to
823 the fact that the search is only partially dependent, in its $O(m \log(n))$ time
824 complexity, on the length m of the pattern while the dominant factor consists
825 of the logarithm of the size n of the data structure.

826 7. Conclusions

827 In this paper we have presented an extensions of a text sampling approach,
828 called Character Distance Sampling, to the case of texts over small alphabets
829 and in the case of offline searching. The first extension was carried out
830 using condensed alphabets in which consecutive groups of q characters are
831 assimilated to a single element of the alphabet, significantly extending its
832 size. The result obtained by this extension was to significantly lower the
833 execution time in the search phase while keeping the space used by the index
834 below the space used by the previous approaches. In our second extension
835 we have proposed a suffix array model built directly on the sampled text in
836 order to decrease the number of candidate occurrences and, consequently, the
837 time required for the response to each single query. This approach contrasts
838 with those currently proposed in the literature which are limited to reducing
839 the number of suffixes taken into account in the original data structure. The
840 experimental results proposed in our extensive experimentation, conducted
841 both on the online and offline scenarios, show how this new proposal offers
842 significantly better performances both in terms of space and in terms of
843 search time. Our future studies will focus in this direction in order to apply
844 sampled string matching to other problems related to text processing.

845 8. Acknowledgement

846 We gratefully acknowledge support from project STORAGE—Università
847 degli Studi di Catania, Piano della Ricerca 2020/2022

848 References

- 849 [1] Ekaterina Benza, Shmuel T. Klein, and Dana Shapira. Smaller compressed suffix arrays†. *Comput. J.*, 64(5):721–730, 2021. doi:10.1093/comjnl/bxaa016.
- 850
851
- 852 [2] Robert S. Boyer and J. Strother Moore. A fast string searching algorithm. *Commun. ACM*, 20(10):762–772, 1977. doi:10.1145/359842.359859.
- 853
854
- 855 [3] Domenico Cantone, Simone Faro, and Emanuele Giaquinta. Adapting boyer-moore-like algorithms for searching huffman encoded texts. *Int. J. Found. Comput. Sci.*, 23(2):343–356, 2012. doi:10.1142/S0129054112400163.
- 856
857
858
- 859 [4] Domenico Cantone, Simone Faro, and Arianna Pavone. Linear and efficient string matching algorithms based on weak factor recognition. *ACM J. Exp. Algorithmics*, 24(1):1.8:1–1.8:20, 2019. doi:10.1145/3301295.
- 860
861
- 862 [5] Francisco Claude, Gonzalo Navarro, Hannu Peltola, Leena Salmela, and Jorma Tarhio. String matching with alphabet sampling. *J. Discrete Algorithms*, 11:37–50, 2012. doi:10.1016/j.jda.2010.09.004.
- 863
864
- 865 [6] Maxime Crochemore, Artur Czumaj, Leszek Gasieniec, Stefan Jarominek, Thierry Lecroq, Wojciech Plandowski, and Wojciech Rytter. Speeding up two string-matching algorithms. *Algorithmica*, 12(4/5):247–267, 1994. doi:10.1007/BF01185427.
- 866
867
868
- 869 [7] Edleno Silva de Moura, Gonzalo Navarro, Nivio Ziviani, and Ricardo A. Baeza-Yates. Fast and flexible word searching on compressed text. *ACM Trans. Inf. Syst.*, 18(2):113–139, 2000. doi:10.1145/348751.348754.
- 870
871
- 872 [8] Simone Faro. A very fast string matching algorithm based on condensed alphabets. In Riccardo Dondi, Guillaume Fertin, and Giancarlo Mauri, editors, *Algorithmic Aspects in Information and Management - 11th International Conference, AAIM 2016, Bergamo, Italy, July 18-20, 2016*,
- 873
874
875

- 876 *Proceedings*, volume 9778 of *Lecture Notes in Computer Science*, pages
877 65–76. Springer, 2016. doi:10.1007/978-3-319-41168-2_6.
- 878 [9] Simone Faro and Thierry Lecroq. The exact online string matching
879 problem: A review of the most recent results. *ACM Comput. Surv.*,
880 45(2):13:1–13:42, 2013. doi:10.1145/2431211.2431212.
- 881 [10] Simone Faro, Thierry Lecroq, Stefano Borzi, Simone Di Mauro, and
882 Alessandro Maggio. The string matching algorithms research tool. In
883 Jan Holub and Jan Zdárek, editors, *Proceedings of the Prague Stringol-
884 ogy Conference 2016, Prague, Czech Republic, August 29-31, 2016*,
885 pages 99–111. Department of Theoretical Computer Science, Faculty
886 of Information Technology, Czech Technical University in Prague, 2016.
887 URL: <http://www.stringology.org/event/2016/p09.html>.
- 888 [11] Simone Faro and Francesco Pio Marino. Reducing time and space
889 in indexed string matching by characters distance text sampling. In
890 Jan Holub and Jan Zdárek, editors, *Prague Stringology Conference
891 2020, Prague, Czech Republic, August 31 - September 2, 2020*, pages
892 148–159. Czech Technical University in Prague, Faculty of Information
893 Technology, Department of Theoretical Computer Science, 2020. URL:
894 <http://www.stringology.org/event/2020/p13.html>.
- 895 [12] Simone Faro, Francesco Pio Marino, and Arianna Pavone. Efficient on-
896 line string matching based on characters distance text sampling. *Algo-
897 rithmica*, 82(11):3390–3412, 2020. doi:10.1007/s00453-020-00732-4.
- 898 [13] Simone Faro, Francesco Pio Marino, and Arianna Pavone. Enhanc-
899 ing characters distance text sampling by condensed alphabets. In
900 Claudio Sacerdoti Coen and Ivano Salvo, editors, *Proceedings of the
901 22nd Italian Conference on Theoretical Computer Science, Bologna,
902 Italy, September 13-15, 2021*, volume 3072 of *CEUR Workshop Pro-
903 ceedings*, pages 1–15. CEUR-WS.org, 2021. URL: [http://ceur-ws.
904 org/Vol-3072/paper1.pdf](http://ceur-ws.org/Vol-3072/paper1.pdf).
- 905 [14] Simone Faro and Stefano Scafiti. The range automaton: An efficient ap-
906 proach to text-searching. In Thierry Lecroq and Svetlana Puzynina, edi-
907 tors, *Combinatorics on Words - 13th International Conference, WORDS
908 2021, Rouen, France, September 13-17, 2021, Proceedings*, volume 12847

- 909 of *Lecture Notes in Computer Science*, pages 91–103. Springer, 2021.
910 doi:10.1007/978-3-030-85088-3_8.
- 911 [15] Paolo Ferragina and Giovanni Manzini. Indexing compressed text. *J.*
912 *ACM*, 52(4):552–581, 2005. doi:10.1145/1082036.1082039.
- 913 [16] Kimmo Fredriksson and Szymon Grabowski. A general compression
914 algorithm that supports fast searching. *Inf. Process. Lett.*, 100(6):226–
915 232, 2006. doi:10.1016/j.ipl.2006.04.020.
- 916 [17] Travis Gagie, Gonzalo Navarro, and Nicola Prezza. Optimal-time text
917 indexing in bwt-runs bounded space. In Artur Czumaj, editor, *Proceed-*
918 *ings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete*
919 *Algorithms, SODA 2018, New Orleans, LA, USA, January 7-10, 2018*,
920 pages 1459–1477. SIAM, 2018. doi:10.1137/1.9781611975031.96.
- 921 [18] Travis Gagie, Gonzalo Navarro, and Nicola Prezza. Fully functional
922 suffix trees and optimal text searching in bwt-runs bounded space. *J.*
923 *ACM*, 67(1):2:1–2:54, 2020. doi:10.1145/3375890.
- 924 [19] Szymon Grabowski and Marcin Raniszewski. Sampled suffix array with
925 minimizers. *Softw. Pract. Exp.*, 47(11):1755–1771, 2017. doi:10.1002/
926 spe.2481.
- 927 [20] R. Nigel Horspool. Practical fast searching in strings. *Softw. Pract.*
928 *Exp.*, 10(6):501–506, 1980. doi:10.1002/spe.4380100608.
- 929 [21] Tomohiro I, Juha Kärkkäinen, and Dominik Kempa. Faster sparse suffix
930 sorting. In Ernst W. Mayr and Natacha Portier, editors, *31st Interna-*
931 *tional Symposium on Theoretical Aspects of Computer Science (STACS*
932 *2014), STACS 2014, March 5-8, 2014, Lyon, France*, volume 25 of
933 *LIPICs*, pages 386–396. Schloss Dagstuhl - Leibniz-Zentrum für Infor-
934 matik, 2014. doi:10.4230/LIPICs.STACS.2014.386.
- 935 [22] Juha Kärkkäinen and Peter Sanders. Simple linear work suffix ar-
936 ray construction. In Jos C. M. Baeten, Jan Karel Lenstra, Joachim
937 Parrow, and Gerhard J. Woeginger, editors, *Automata, Languages and*
938 *Programming, 30th International Colloquium, ICALP 2003, Eindhoven,*
939 *The Netherlands, June 30 - July 4, 2003. Proceedings*, volume 2719
940 of *Lecture Notes in Computer Science*, pages 943–955. Springer, 2003.
941 doi:10.1007/3-540-45061-0_73.

- 942 [23] Shmuel T. Klein and Dana Shapira. A new compression method for
943 compressed matching. In *Data Compression Conference, DCC 2000,*
944 *Snowbird, Utah, USA, March 28-30, 2000*, pages 400–409. IEEE Com-
945 puter Society, 2000. doi:10.1109/DCC.2000.838180.
- 946 [24] Donald E. Knuth, James H. Morris Jr., and Vaughan R. Pratt. Fast
947 pattern matching in strings. *SIAM J. Comput.*, 6(2):323–350, 1977.
948 doi:10.1137/0206024.
- 949 [25] Tomasz Marek Kowalski, Szymon Grabowski, and Kimmo Fredriksson.
950 Suffix arrays with a twist. *Comput. Informatics*, 38(3):555–574, 2019.
951 URL: [http://www.cai.sk/ojs/index.php/cai/article/view/2019_](http://www.cai.sk/ojs/index.php/cai/article/view/2019_3_555)
952 [3_555](http://www.cai.sk/ojs/index.php/cai/article/view/2019_3_555).
- 953 [26] Alan Kuhnle, Taher Mun, Christina Boucher, Travis Gagie, Ben Lang-
954 mead, and Giovanni Manzini. Efficient construction of a complete index
955 for pan-genomics read alignment. *J. Comput. Biol.*, 27(4):500–513, 2020.
956 doi:10.1089/cmb.2019.0309.
- 957 [27] Zhize Li, Jian Li, and Hongwei Huo. Optimal in-place suffix sorting. In
958 Travis Gagie, Alistair Moffat, Gonzalo Navarro, and Ernesto Cuadros-
959 Vargas, editors, *String Processing and Information Retrieval - 25th In-*
960 *ternational Symposium, SPIRE 2018, Lima, Peru, October 9-11, 2018,*
961 *Proceedings*, volume 11147 of *Lecture Notes in Computer Science*, pages
962 268–284. Springer, 2018. doi:10.1007/978-3-030-00479-8_22.
- 963 [28] Veli Mäkinen and Gonzalo Navarro. Succinct suffix arrays based on
964 run-length encoding. In Alberto Apostolico, Maxime Crochemore, and
965 Kunsoo Park, editors, *Combinatorial Pattern Matching*, pages 45–56,
966 Berlin, Heidelberg, 2005. Springer Berlin Heidelberg.
- 967 [29] Udi Manber. A text compression scheme that allows fast searching
968 directly in the compressed file. *ACM Trans. Inf. Syst.*, 15(2):124–136,
969 1997. doi:10.1145/248625.248639.
- 970 [30] Udi Manber and Gene Myers. Suffix arrays: A new method for on-
971 line string searches. *SIAM J. Comput.*, 22(5):935–948, oct 1993. doi:
972 10.1137/0222058.
- 973 [31] Gonzalo Navarro. *Compact Data Structures - A practical approach*. 2016.

- 974 [32] Gonzalo Navarro. The compression power of the BWT: technical per-
975 spective. *Commun. ACM*, 65(6):90, 2022. doi:10.1145/3531443.
- 976 [33] Gonzalo Navarro and Jorma Tarhio. Lzgrep: a boyer-moore string
977 matching tool for ziv-lempel compressed text. *Softw. Pract. Exp.*,
978 35(12):1107–1130, 2005. doi:10.1002/spe.663.
- 979 [34] Simon J. Puglisi and Bella Zhukova. Relative lempel-ziv compression of
980 suffix arrays. In Christina Boucher and Sharma V. Thankachan, editors,
981 *String Processing and Information Retrieval - 27th International Sym-*
982 *posium, SPIRE 2020, Orlando, FL, USA, October 13-15, 2020, Proceed-*
983 *ings*, volume 12303 of *Lecture Notes in Computer Science*, pages 89–96.
984 Springer, 2020. doi:10.1007/978-3-030-59212-7_7.
- 985 [35] Simon J. Puglisi and Bella Zhukova. Smaller rlz-compressed suffix ar-
986 rays. In Ali Bilgin, Michael W. Marcellin, Joan Serra-Sagristà, and
987 James A. Storer, editors, *31st Data Compression Conference, DCC*
988 *2021, Snowbird, UT, USA, March 23-26, 2021*, pages 213–222. IEEE,
989 2021. doi:10.1109/DCC50243.2021.00029.
- 990 [36] Yusuke Shibata, Takuya Kida, Shuichi Fukamachi, Masayuki Takeda,
991 Ayumi Shinohara, Takeshi Shinohara, and Setsuo Arikawa. Speeding
992 up pattern matching by text compression. In Gian Carlo Bongiovanni,
993 Giorgio Gambosi, and Rossella Petreschi, editors, *Algorithms and Com-*
994 *plexity, 4th Italian Conference, CIAC 2000, Rome, Italy, March 2000,*
995 *Proceedings*, volume 1767 of *Lecture Notes in Computer Science*, pages
996 306–315. Springer, 2000. doi:10.1007/3-540-46521-9_25.
- 997 [37] Uzi Vishkin. Deterministic sampling - A new technique for fast pattern
998 matching. *SIAM J. Comput.*, 20(1):22–40, 1991. doi:10.1137/0220002.
- 999 [38] Andrew Chi-Chih Yao. The complexity of pattern matching for a ran-
1000 dom string. *SIAM J. Comput.*, 8(3):368–387, 1979. doi:10.1137/
1001 0208029.