

Adaptive Learning of Compressible Strings

Gabriele Fici^a, Nicola Prezza^b, Rossano Venturini^c

^a*Dipartimento di Matematica e Informatica, Università degli Studi di Palermo, Palermo, Italy*
gabriele.fici@unipa.it

^b*Dipartimento di Scienze Ambientali, Informatica e Statistica, Università Ca' Foscari, Venezia, Italy*
nicola.prezza@unive.it

^c*Dipartimento di Informatica, Università di Pisa, Pisa, Italy*
rossano.venturini@unipi.it

Abstract

Suppose an oracle knows a string S that is unknown to us and that we want to determine. The oracle can answer queries of the form “Is s a substring of S ?”. In 1995, Skiena and Sundaram showed that, in the worst case, any algorithm needs to ask the oracle $\sigma n/4 - O(n)$ queries in order to be able to reconstruct the hidden string, where σ is the size of the alphabet of S and n its length, and gave an algorithm that spends $(\sigma - 1)n + O(\sigma\sqrt{n})$ queries to reconstruct S .

The main contribution of our paper is to improve the above upper-bound in the context where the string is compressible. We first present a universal algorithm that, given a (computable) compressor that compresses the string to τ bits, performs $q = O(\tau)$ substring queries; this algorithm, however, runs in exponential time. For this reason, the second part of the paper focuses on more time-efficient algorithms whose number of queries is bounded by specific compressibility measures. We first show that any string of length n over an integer alphabet of size σ with rle runs can be reconstructed with $q = O(\text{rle}(\sigma + \log \frac{n}{\text{rle}}))$ substring queries in linear time and space. We then present an algorithm that spends $q \in O(\sigma g \log n)$ substring queries and runs in $O(n(\log n + \log \sigma) + q)$ time using linear space, where g is the size of a smallest straight-line program generating the string.

Keywords: String reconstruction, String learning, Adaptive learning, Kolmogorov complexity, String Compression, Lempel-Ziv, Centroid decomposition, Suffix tree.

1. Introduction

String reconstruction (or learning) from substrings queries is a well-established problem that has natural applications in many areas, including bioinformatics, data compression, security, etc. (see, for example, [1, 2, 3, 4]).

In a more general setting, one is interested in understanding whether and how it is possible to reconstruct an unknown target string S from some piece of information about S . This information can be, for example, a collection of substrings (e.g., the classical NP-hard Shortest

*Supported by MIUR PRIN 2017 Project 2017K7XPAN and Università di Pisa Project PRA.2020-2021.26.

8 Superstring Problem), or substring compositions ([5]), or subwords ([6]), of S . Furthermore,
9 the problem can be viewed from different angles, e.g., combinatorial, computational, algo-
10 rithmic, information theoretical.

11 In this paper, we deal with the problem of reconstructing a string from information about
12 its substrings. Apart from the classical static model for the reconstruction (exact or with
13 uncertainty), many different models have been introduced in the literature for the string
14 reconstruction problem, including the one we consider in this paper, and which has been
15 presented in 1995 by Skiena and Sundaram [3]. In this model, one can ask an oracle, which
16 knows the target string S , queries of the form “Is s a substring of S ?” and is interested
17 in designing an *adaptive algorithm* minimizing the number of such queries. In this setting,
18 *adaptive* means that the algorithm may reuse the information resulting from previous queries
19 in order to decide which queries to ask next. It is worth mentioning that, with the same
20 model, other query complexities have been investigated very recently by Amir et al. [7].

21 A trivial information-theoretic argument implies a worst-case lower bound of $n \log \sigma$
22 queries, where σ is the size of the alphabet of S . Skiena and Sundaram [3] improved this
23 bound and showed that $\sigma n/4 - O(n)$ queries are necessary to reconstruct S in the worst
24 case. This remains true even if the oracle returns for each substring query the number of its
25 occurrences in S . In the same paper, they gave an algorithm for the reconstruction which
26 spends at most $(\sigma - 1)n + O(\sigma\sqrt{n})$ queries, thus asymptotically matching the lower bound.
27 They also gave an algorithm that spends at most $(\sigma - 1)n + 2 \log n + O(\sigma)$ queries if the length
28 n of S is known. Iwama et al. gave an algorithm for binary strings that spends $n + O(1)$
29 queries *on average* [8]. Amir et al. [7] recently proved that if the string has period $p > 0$,
30 then it can be reconstructed using $O(\sigma p + \lg n)$ substring queries, even if both n and p are
31 unknown.

32 We stress out that these bounds hold in the *adaptive* case: answers to previous queries
33 can be used to decide the next query. As shown by Skiena et al. [2], the non-adaptive model
34 is much harder: if the algorithm has to reconstruct all substrings (of the unknown string) of
35 length $k \leq n$ after a pre-determined batch of queries, then $\sigma^{k/2}/k$ queries are needed to solve
36 the problem. Tsur [4] explored this model more in detail, providing bounds as a function of
37 the fraction $(1 - \epsilon)$ (with $0 \leq \epsilon \leq 1$) of substrings of length k that have to be reconstructed:
38 in this case, $\Omega(\epsilon^{-1/2}k^2)$ non-adaptive queries are sufficient and necessary.

39 1.1. A novel approach to the problem

40 While the aforementioned papers tackled the problem in the worst-case, the minimum
41 number of queries needed to reconstruct a string may be significantly smaller than the worst-
42 case on particular instances. For example, consider a string of the form a^n , where a is a
43 single letter. An algorithm could first try to find out if the string is of this form by issu-
44 ing $O(\log n) + 2\sigma$ queries and, only if the check fails, proceed with Skiena and Sundaram’s
45 algorithm [3]. Observe that the resulting algorithm optimizes for a particular class of *highly-*
46 *compressible* strings. In fact, in this paper we show that this reasoning continues to hold for
47 *any compressor*. Our first result is a universal algorithm that, given as input a computable
48 compressor \mathcal{C} , performs the reconstruction asking a number of queries that is proportional
49 to the bit-size $|\mathcal{C}(S)|$ of the string compressed by \mathcal{C} . We complement this result by showing

50 that any deterministic adaptive algorithm for reconstructing a string yields a string com-
 51 pressor. Together, these results imply the equivalence between the string reconstruction and
 52 compression problems.

53 Motivated by the fact that our universal algorithm performs an exponential number of
 54 calls to \mathcal{C} , we then focus on optimizing the running time and the space usage for commonly
 55 used compressors, including run-length encoding, Lempel-Ziv factorization and context-free
 56 grammars.

57 In measuring the efficiency of an algorithm, we assume that any query can be submitted
 58 to the oracle in constant time and space regardless of the length of the queried substring. The
 59 reason for this assumption is that the implementation of the oracle strongly depends on the
 60 application. For example, if the application admits a collaborative oracle, there are several
 61 possible approaches to achieve constant query time, e.g., using hashing. Moreover, one could
 62 also assume that the oracle knows the reconstruction strategy and therefore it could run the
 63 reconstruction algorithm itself, that is, we do not even need to transmit the next substring
 64 query because the oracle already knows the next query it has to answer.

65 1.2. Preliminary definitions

66 Let S be a binary string. A *compressor* is an injective computable function $\mathcal{C} : \{0, 1\}^+ \rightarrow$
 67 $\{0, 1\}^+$ that converts any $S \in \{0, 1\}^+$ into a reversible representation $\mathcal{C}(S)$ of size $|\mathcal{C}(S)|$
 68 bits. We require also the inverse function \mathcal{C}^{-1} (i.e. the function such that $\mathcal{C}^{-1}(\mathcal{C}(S)) = S$)
 69 to be computable; this function is the *decompressor* associated with \mathcal{C} . Informally speaking,
 70 a function \mathcal{C} qualifies as a good compressor if $|\mathcal{C}(S)| \ll |S|$ on particular string families (for
 71 example, repetitive strings), and $|\mathcal{C}(S)| \in O(|S|)$ for all other strings outside this family.

72 A popular compressor is the LZ77 factorization of S . The Lempel-Ziv 1977 (LZ77) algo-
 73 rithm [9] parses a string S into a sequence of z phrases, where each new phrase is either a
 74 fresh character or the longest string that also occurs starting from a position strictly smaller
 75 than the phrase start position. The bit-size of the LZ77 factorization of S is $\Theta(z \log n)$. For
 76 example, the LZ77 factorization of the string *abbabba* is *a|b|b|abba*. In this example, the string
 77 is factored into $z = 4$ phrases. A more restricted version of LZ77 does not allow overlaps
 78 between a phrase and its source. We denote this version as *LZ77 without overlaps* and with
 79 z_{no} the number of generated phrases. Between those two measures, it holds $z_{no} \in O(z \log n)$
 80 [10]. Clearly, also $z \leq z_{no}$ always holds. This version factorizes the above string as *a|b|b|abb|a*,
 81 with $z_{no} = 5$.

82 Another common measure of compressibility is the number *rle* of equal-letter runs in S ,
 83 that is, the number of maximal unary substrings of S . This measure is not as strong as z ; in
 84 fact, it is easy to see that $z_{no} \leq \text{rle}$.

85 In this work we also consider the size (number of nonterminals) g of the smallest straight-
 86 line program (SLP) producing (only) S . SLPs are particular cases of acyclic context-free
 87 grammars composed of rules of the kind $A \rightarrow BC$, where A is a nonterminal and B, C are
 88 either nonterminals or terminals.

89 The known relations between g and z_{no} are $g \in O(z_{no} \log(n/z_{no}))$ and $z_{no} \leq g$. See
 90 Navarro's recent survey [10] for more details on these and several other relations between
 91 string complexity measures.

92 **2. Universal string reconstruction**

93 In this section we present an algorithm that, given a compressor C , reconstructs any
 94 string S with $O(|C(S)|)$ queries to the oracle. We furthermore prove a dual result: any
 95 reconstruction algorithm performing $\chi(S)$ queries yields a compression algorithm (with asso-
 96 ciated decompressor) that compresses the string to $\chi(S)$ bits. These findings show that the
 97 string reconstruction problem essentially coincides with the string compression problem. For
 98 simplicity, in this section we restrict our attention to binary alphabets only.

99 We start with a lemma of Skiena and Sundaram [3] stating that any set M of binary
 100 strings admits a string that is a substring of a constant fraction of the members of M .
 101 Letting M be a set of strings, we let $M(S)$ denote the subset of M whose elements contain
 102 S as a substring.

103 **Lemma 1.** (*[3, Lem. 12]*) *Let $M \subseteq \{0, 1\}^n$ be a set of binary strings, each of length n .*
 104 *Then, there exists a string S such that $\frac{1}{5}|M| \leq |M(S)| \leq \frac{4}{5}|M|$.*

105 Lemma 1 can be turned into a universal algorithm for determining the substring queries
 106 to be asked to the oracle as a function of any given compressor.

107 **Lemma 2.** *Let \mathcal{C} be a compressor, and let $S \in \{0, 1\}^n$ be an unknown binary string of known
 108 length n . Then, there is an algorithm that reconstructs S using $O(|C(S)|)$ substring queries.*

109 *Proof.* Let $M_k = \{S \in \{0, 1\}^n : |C(S)| \leq k\}$ be the set of strings of length n compressed to
 110 at most k bits by \mathcal{C} . Note that $|M_k| \leq 2^{k+1} - 2$, since \mathcal{C} is injective and there are no more
 111 than $2^{k+1} - 2$ binary strings (compressed representations) of length at most k . Assuming
 112 we know the value of $\tau = |C(S)|$ (later we remove this assumption), it is easy to design
 113 an (exponential-time) algorithm that builds M_τ : simply apply \mathcal{C} to all strings of length n ,
 114 keeping only those such that $|C(S)| \leq \tau$. By definition of τ , note that $S \in M_\tau$. Then,
 115 by applying recursively Lemma 1 starting from the set M_τ , we end up selecting S from
 116 this set. Each recursive iteration yields a string that we use to perform a substring query
 117 on S , thereby reducing the number of candidates by a factor $4/5$ in the worst case. After
 118 $O(\log |M_\tau|) = O(\tau)$ iterations (i.e., substring queries on S), we discover which element of M_τ
 119 corresponds to S . To conclude, we can remove the assumption that we know τ . To achieve
 120 this goal it is sufficient to run an exponential search on the above strategy, i.e., run it on $M_{\tau'}$
 121 for $\tau' = 1, 2, 4, \dots, 2^{\lceil \log_2 \tau \rceil}$. The last iteration will reveal S , after a total of $O(\tau)$ substring
 122 queries on S . \square

123 We finally prove the following lemma.

124 **Lemma 3.** *Let A be a deterministic adaptive algorithm that reconstructs any string S by
 125 asking $\chi(S)$ queries to the oracle, for some (computable) function $\chi(S)$. Then, there exists
 126 a compressor \mathcal{C} (and an associated decompressor \mathcal{C}^{-1}) such that $|C(S)| = \chi(S)$.*

127 *Proof.* It is straightforward to turn A into a compressor \mathcal{C} : the compressed representation
 128 $\mathcal{C}(S)$ of S is the binary string of length $\chi(S)$ formed by the $\chi(S)$ answers received by the
 129 oracle while reconstructing S . The $\chi(S)$ answers can be computed by any pattern matching

130 algorithm testing membership of the substrings queried by A in the substring closure of S .
 131 Similarly, A itself can be turned into a decompressor: by definition of A , the $\chi(S)$ answers
 132 of the oracle (i.e. the compressed file representation) are sufficient to reconstruct A . \square

133 While the above results establish an asymptotic equivalence between the string reconstruc-
 134 tion and compression problems, they do not yield time-efficient algorithms for reconstructing
 135 a string in time proportional to its compressed size. In the next section we tackle this problem
 136 by focusing on particular string compressors.

137 3. Feasible algorithms for the reconstruction

138 Let S be a string of length n over an integer alphabet $\Sigma = [1, \dots, \sigma]$. A trivial algorithm
 139 for reconstructing S with $\sigma(n + 1)$ substring queries is the following [3]: We make queries of
 140 single character substrings, so that after at most σ queries a new character of S is determined.
 141 Let s be a known substring of S . In general, we can increase the length of this known substring
 142 by one character by querying on the strings $s\sigma_i$, for every character σ_i . At least one of these
 143 queries must be a substring of S , unless s is a suffix of S that has no other occurrences in S .
 144 When s can no longer be extended to the right, i.e., s is a suffix of S not appearing elsewhere
 145 in S , we can continue the process by prepending characters to the known substring s , until
 146 it can no longer be extended to the left, and the string S is then reconstructed.

147 This algorithm is optimal up to constant factors due to the following lower bound [3].

148 **Theorem 4.** ([3, Thm. 8]) *In the worst case, $\frac{\sigma n}{4} - O(n)$ substring queries are necessary to*
 149 *reconstruct a string of length n .*

150 In the rest of this section, we will provide algorithms for reconstructing the string S whose
 151 efficiency is measured towards commonly used measures of compression for strings.

152 Let us first show an easy result for the size rle of the run-length encoding of S , i.e., rle is
 153 the number of runs (maximal repetitions of the same character) in S . We show that S can
 154 be reconstructed with $O(\text{rle}(\sigma + \log \frac{n}{\text{rle}}))$ queries. The reconstruction is done in rle steps. Let
 155 \hat{S}_{i-1} be the substring reconstructed so far. In the i th step, we first identify the character c
 156 that follows \hat{S}_{i-1} in S . This is done by querying $\hat{S}_{i-1} \cdot c$ for any $c \in \Sigma$. Once we know c , we
 157 need to identify the length of the run of c , i.e., the maximal value r_i such that $\hat{S}_{i-1} \cdot c^{r_i}$ is
 158 a substring of S . This can be done with an exponential search on r_i , which takes $\Theta(\log r_i)$
 159 queries. When at some step j , \hat{S}_j cannot be extended further, we continue the process by
 160 prepending (runs of) characters to the known substring.

161 The overall number of queries is $q = O(\sum_{i=1}^{\text{rle}} (\sigma + \log r_i))$. This is in $O(\text{rle}(\sigma + \log \frac{n}{\text{rle}}))$
 162 queries because the sum of the terms $\log r_i$ is maximized when every r_i is in $\Theta(\frac{n}{\text{rle}})$.

163 **Theorem 5.** *Any string of length n with rle runs can be reconstructed with $q = O(\text{rle}(\sigma +$
 164 $\log \frac{n}{\text{rle}}))$ substring queries in $O(q)$ time and $O(\text{rle})$ space.*

165 Note that, accordingly to Theorem 4, this result is optimal up to a constant factor for
 166 sufficiently large σ .

167 Our next aim is to give algorithms whose complexity grows as a function of the size of
 168 the LZ77 parsing of S . We let z_{no} denote the number of phrases of the LZ77 parsing when
 169 the parse does not allow overlapping phrases (both settings are commonly considered in the
 170 literature).

171 We can use Theorem 4 to prove a lower bound on $\chi(S)$ in terms of z_{no} .

172 **Theorem 6.** *In the worst case, $\Omega(\sigma z_{no} \log_{\sigma} n)$ substring queries are necessary to reconstruct
 173 a string of length n .*

174 *Proof.* It is well known that for any string $z_{no} = O(\frac{n}{\log_{\sigma} n})$. The theorem follows by combining
 175 this fact with Theorem 4. \square

176 We are not required to know the length n of S , but we assume to know its alphabet
 177 $\Sigma = [1, \dots, \sigma]$. If instead also the alphabet is unknown, we need $O(\log \sigma)$ queries to identify
 178 the largest character in S . This is done by performing an exponential search to identify the
 179 largest character occurring in S . Notice that this is correct only if all the characters in Σ
 180 occur in S (in particular, $\sigma \leq n$), which we assume as hypothesis.

181 Our goal is to prove the following theorem.

182 **Theorem 7.** *Let S be a string of length n over the alphabet $\Sigma = [1, \dots, \sigma]$. There exists
 183 an algorithm that reconstructs S with $q = O(\sigma z_{no} \log(n/z_{no}) \log n)$ substring queries to the
 184 oracle. The algorithm runs in $O(n(\log n + \log \sigma) + q)$ time using linear space.*

185 Note that this result is optimal up to a factor $O(\log \frac{n}{z_{no}} \log \sigma)$ by Theorem 6.

186 In the next subsection we review a technique to solve pattern matching queries on a text
 187 which exploits the centroid decomposition of the suffix tree of a string. This will allow us to
 188 give an efficient algorithm for reconstructing the suffix tree of S , from which S is therefore
 189 determined.

190 *Pattern matching with the centroid decomposition.* This technique has been introduced by
 191 Naor [11] and has found applications, for example, in designing cache-oblivious string B-
 192 trees [12, 13] or randomized pattern matching [14] on a dictionary of strings.

193 The *centroid decomposition* of a tree \mathcal{T} (also known as *separator decomposition*) is a
 194 popular and powerful technique to obtain a tree $\mathcal{T}_{\mathcal{C}}$ of logarithmic height. The decomposition
 195 is based on a theorem proved by Jordan in 1869 [15].

196 **Lemma 8.** *Any tree \mathcal{T} of n nodes has at least a node, called centroid, whose removal leaves
 197 connected components of size at most $n/2$.*

198 The centroid decomposition is defined recursively. Given \mathcal{T} , we identify a centroid node
 199 u , which is chosen to be the root of the new rooted tree $\mathcal{T}_{\mathcal{C}}$. Then, we remove u from \mathcal{T} and
 200 recurse on each connected component to get u 's subtrees in $\mathcal{T}_{\mathcal{C}}$. The children of u in $\mathcal{T}_{\mathcal{C}}$ are
 201 the roots of the centroid decompositions of these components. Let us use $\text{children}_{\mathcal{T}_{\mathcal{C}}}(u)$ to
 202 denote the set of children of u in $\mathcal{T}_{\mathcal{C}}$. As we have a (possibly empty) component for u 's parent
 203 and children in \mathcal{T} , the outdegree of u in $\mathcal{T}_{\mathcal{C}}$, i.e., $|\text{children}_{\mathcal{T}_{\mathcal{C}}}(u)|$, is at most the outdegree of

204 u in \mathcal{T} plus one. The resulting decomposition is a new tree \mathcal{T}_C on the same nodes, whose
205 height is $\Theta(\log n)$.

206 A folklore algorithm computes the centroid decomposition in $\Theta(n \log n)$ time as follows.
207 We first observe that a centroid node of \mathcal{T} can be easily identified in linear time. Indeed, we
208 can arbitrary choose a root in \mathcal{T} and visit the tree to compute the size of each subtree. Then,
209 we start from the root and move to the largest subtree until we reach a node whose subtrees
210 have size at most $n/2$. This node is a centroid of the tree. The centroid decomposition is
211 computed by repeating the above algorithm recursively in each component. It easily follows
212 that the decomposition of the tree can be computed in $\Theta(n \log n)$ time. However, there exist
213 construction algorithms to compute the decomposition in linear time [16, 17].

214 In the following, we will use the centroid decomposition \mathcal{ST}_C of the suffix tree \mathcal{ST} of a
215 string S . Given a node u in \mathcal{ST} , we use $\text{locus}(u)$ to denote its locus, i.e., the string obtained
216 by concatenating the sequence of labels encountered along the path from the root to u .

217 Assume we are given a pattern $P[1, p]$ and our goal is to find the suffix of S that shares
218 the longest common prefix with P . This problem can be easily solved with the suffix tree \mathcal{ST}
219 of S with the following two-phase strategy: We first identify the highest node u^* in \mathcal{ST} such
220 that $\text{locus}(u^*)$ shares the longest common prefix with P . Then, we try to extend the match
221 by comparing the remaining characters of P with the characters on the edge between u^* and
222 one of its children, i.e., the child where the label starts with the character $P[|\text{locus}(u^*)| + 1]$.

223 We can perform the same search for u^* on the centroid decomposition \mathcal{ST}_C of the suffix
224 tree of S . The search is done by traversing a root-to-node path of $O(\log n)$ nodes. We start
225 from the root of \mathcal{ST}_C and we move down to the leaves. For every node u we visit, we compare
226 $\text{locus}(u)$ with P and decide in which of its children we have to continue the search. As the
227 target node u^* is guaranteed to be visited, we simply take track of the visited node sharing
228 the longest common prefix with P . Based on the result of comparing $\text{locus}(u)$ and P , there
229 are the following cases:

- 230 • If $\text{locus}(u)$ equals P , then u is our target node u^* and we conclude.
- 231 • If $\text{locus}(u)$ is not a prefix of P , we continue to search on the child of u which corresponds
232 to the connected component containing the parent of u in the suffix tree. If such a node
233 does not exist, we conclude.
- 234 • If $\text{locus}(u)$ is a prefix of P , we continue the search on the connected component con-
235 taining one of the children of u in the suffix tree. The child is the node v such that
236 the first character of the edge between u and v equals the character $P[|\text{locus}(u)| + 1]$.
237 This is exactly the node v that a normal search on the suffix tree would visit next,
238 once the search reaches u . Notice that in general v is not a child of u in the centroid
239 decomposition. If such a node does not exist, we finish the visit.

240 Let us now show how to use the above algorithm to reconstruct an unknown string S .

241 *Solution with prefix queries to the oracle.* We first describe our algorithm for querying the
242 oracle in an easier setting. Instead of answering substring queries, the oracle answers *prefix*

243 *queries*: given a string P , the oracle tells us whether P is a prefix of the unknown string S .
 244 This model is stronger because it allows us to remain anchored to the beginning of S while
 245 reconstructing it.¹ A direct consequence is that the algorithm is easier and faster.

246 We now describe how to reconstruct a string S with $\Theta(\sigma z_{no} \log n)$ prefix queries to the
 247 oracle.

248 Our algorithm works in steps. In the i -th step it reconstructs the i -th LZ77 phrase Z_i .
 249 Once Z_i is reconstructed, the algorithm knows the string \hat{S}_i , which is the concatenation of all
 250 the phrases reconstructed so far. Observe that Z_i is the longest substring of \hat{S}_{i-1} such that
 251 the prefix query $Q_i = \hat{S}_{i-1} \cdot Z_i$ is answered affirmatively.

252 The phrase Z_i is identified with $O(\sigma(\log |\hat{S}_{i-1}| + 1))$ prefix queries as follows. Assume
 253 we have the suffix tree \mathcal{ST} of \hat{S}_{i-1} and its centroid decomposition $\mathcal{ST}_{\mathcal{C}}$. Our first goal is to
 254 identify the lowest node u^* in \mathcal{ST} such that $\hat{S}_{i-1} \cdot \text{locus}(u^*)$ is a prefix of S . This can be
 255 done by performing a search for the unknown pattern $P = \text{locus}(u^*)$ on $\mathcal{ST}_{\mathcal{C}}$. Even if u^* is
 256 unknown, the search can be performed correctly. Indeed, observe that u^* and all its ancestors
 257 in \mathcal{ST} are the only nodes u such that $\hat{S}_{i-1} \cdot \text{locus}(u)$ is a prefix of S . Thus, we perform the
 258 search on the centroid tree by binary searching for u^* on root-to- u^* path. The cost of the
 259 search is $O(\sigma(\log |\hat{S}_{i-1}| + 1))$ prefix queries. Indeed, we need to visit $O(\log |\hat{S}_{i-1}| + 1)$ nodes
 260 of $\mathcal{ST}_{\mathcal{C}}$ to identify u^* . For each visited node u , we need a query to check if $\hat{S}_{i-1} \cdot \text{locus}(u)$ is
 261 a prefix of S . If this is the case, at most σ queries of the form $\hat{S}_{i-1} \cdot \text{locus}(u) \cdot c$, with $c \in \Sigma$,
 262 are needed to know in which child of u we have to continue our search. Otherwise, we move
 263 to the component containing the parent of u , if any.

264 Once we know u^* , we have to extend $\text{locus}(u^*)$ to match Z_i . Indeed, Z_i may end up in the
 265 middle of the edge from u^* to one of its children, say v , in \mathcal{ST} . This step can be easily done
 266 with $O(\sigma + \log |Z_i|)$ queries. First, we use $O(\sigma)$ queries to identify the child v of u , then we
 267 perform an exponential search on the length of the edge label.

268 We conclude by proving that the reconstruction of S takes $O(n \log n + n \log \sigma)$ time. A
 269 trivial implementation of our algorithm consists in rebuilding at each step i the suffix tree of
 270 string \hat{S}_i and its centroid decomposition from scratch. This takes quadratic time.

271 A faster algorithm is the following: First, we observe that, as the string is reconstructed
 272 from left to right, we can use the Ukkonen's construction of the suffix tree [18]. This con-
 273 struction builds the suffix tree in $O(n \log \sigma)$ time and linear space. It is an online algorithm
 274 that processes the string from left to right, hence it allows us to build the suffix tree of the
 275 prefix of the string that we have already reconstructed.

276 The centroid decomposition of the suffix tree is instead kept updated dynamically. Brodal
 277 et al. [16] showed how to keep an approximation of the centroid decomposition of a tree sub-
 278 ject to insertions of new nodes in $O(\log n)$ amortized time per insertion. The decomposition is
 279 approximated in the sense that each selected centroid node splits the tree in connected com-
 280 ponents having a fraction $\frac{1}{2} + \epsilon$ of the overall tree dimension, for any $0 < \epsilon < \frac{1}{4}$. The height
 281 of the $\mathcal{T}_{\mathcal{C}}$ is still $O(\log n)$, thus this approximated decomposition suffices for our purposes.

¹An oracle for prefix queries can be obtained from an oracle for substring queries if we assume that S begins with a special character $\$$ not belonging to Σ .

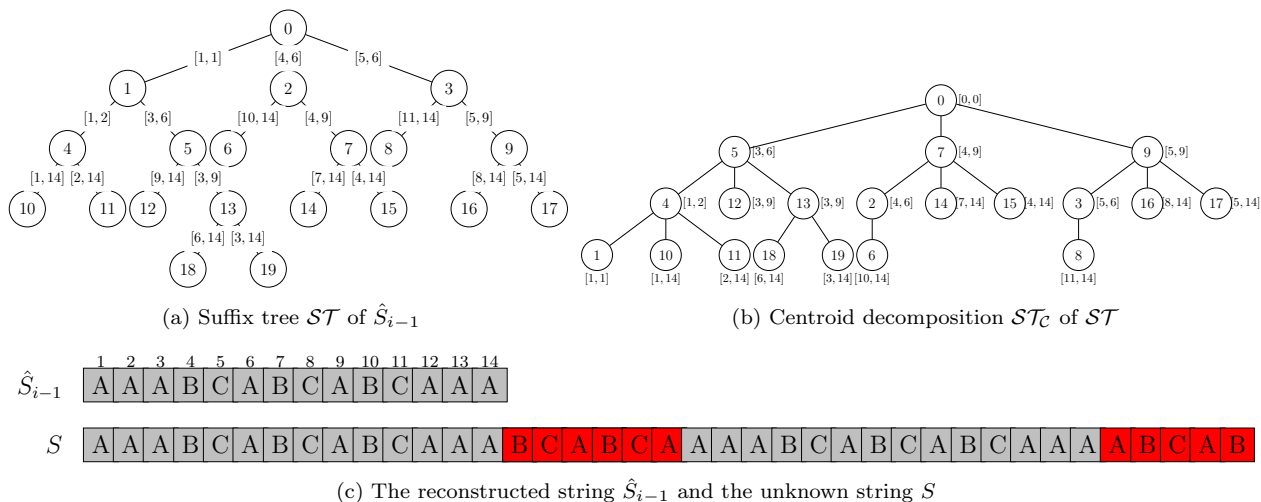


Figure 1: A running example

282 *Solution with substring queries to the oracle.* The string S can be reconstructed with sub-
 283 string queries using an easy variant of the above algorithm. The algorithm reconstructs (a
 284 portion of) the string exactly as described above, until the string reconstructed so far, say
 285 \hat{S}_i , cannot be further extended to the right, hence we know that \hat{S}_i is a suffix of S . Then, we
 286 start extending \hat{S}_i from its beginning, proceeding backwards (that is, prepending characters
 287 to \hat{S}_i). More formally, our strategy first queries forward strings and builds the suffix tree
 288 and the LZ77 factorization of some suffix $S[i, n]$ of the string. Then, we build the suffix tree
 289 of $\overleftarrow{S[i, n]}$ and proceed backwards, building the LZ77 factorization of the remaining portion
 290 $\overleftarrow{S[1, i-1]}$. Since the size g of the smallest grammar is invariant under reversals and upper-
 291 bounds the number of Lempel-Ziv phrases, in both phases we generate at most $O(g)$ phrases.
 292 The following theorem is therefore immediate.

293 **Theorem 9.** *Let S be a string of length n over the alphabet $\Sigma = [1, \dots, \sigma]$. There exists*
 294 *an algorithm that reconstructs S with $q = O(\sigma g \log n)$ substring queries to the oracle. The*
 295 *algorithm runs in $O(n(\log n + \log \sigma) + q)$ time using linear space.*

296 Finally, Theorem 7 follows from the well-known bound $g \in O(z_{no} \log(n/z_{no}))$ (see also
 297 Navarro [10]).

298 *Running example.* Suppose we have already reconstructed the string $\hat{S}_{i-1} = \text{AAABCABCABCAAA}$.
 299 This is a substring of the unknown string S shown in Figure 1c. There are two occurrences
 300 of \hat{S}_{i-1} in S and the red cells highlight the characters that we still need to learn. The suffix
 301 tree \mathcal{ST} of \hat{S}_{i-1} (see Figure 1a) has been built online with Ukkonen's algorithm and it will be
 302 updated as soon as we learn more characters. For this reason we do not append any special
 303 character at the end of \hat{S}_{i-1} . Thus, there may exist suffixes of \hat{S}_{i-1} which do not have their
 304 leaves in \mathcal{ST} because they are proper prefixes of some another suffix. In our example this
 305 happens to the last three suffixes A, AA and AAA which are proper prefixes of the whole string

306 \hat{S}_{i-1} . Nodes of \mathcal{ST} are numbered (in our example levelwise just for convenience) to map the
 307 corresponding node in the centroid decomposition \mathcal{ST}_C (see Figure 1b).

308 The label on the edge from node u to its child v reports the interval $[i, j]$ of positions on
 309 string \hat{S}_{i-1} representing the locus of node v . For example, the label on the edge from node
 310 5 to node 13 is $[3, 9]$ and, thus, $\text{locus}(13) = \text{ABCABCA}$. In the centroid decomposition \mathcal{ST}_C ,
 311 each node u is labeled with the interval of positions of $\text{locus}(u)$. The label of node 13 is $[3, 9]$
 312 because $\text{locus}(13) = \text{ABCABCA}$. In the centroid tree, the leftmost child of any node u is the
 313 centroid decomposition of the connected component containing the parent of u (if any) while
 314 the other children are the centroid decompositions of the subtrees rooted at the children of
 315 u in \mathcal{ST} (if any). Note that for any child v of node u but, possibly, the leftmost one, we have
 316 that $\text{locus}(u)$ is a prefix of $\text{locus}(v)$.

317 We start from the root of \mathcal{ST}_C (node 0) which in our example, by coincidence, corresponds
 318 to the root of \mathcal{ST} . We query the oracle for the substring $\hat{S}_{i-1} \cdot \text{locus}(0)$. As $\text{locus}(0)$ is the
 319 empty string, the oracle's answer will be positive. The algorithm continues on a child of node
 320 0.

321 For any child v of u , let be c_v the character such that $\text{locus}(u) \cdot c_v$ is a prefix of $\text{locus}(v)$,
 322 i.e., $c_v = \text{locus}(v)[|\text{locus}(u)| + 1]$. We process the children of node 0 and continue on a node
 323 v such that the query for the substring $\hat{S}_{i-1} \cdot \text{locus}(0) \cdot c_v$ is successful. There may be several
 324 such nodes v . For example, we could continue on both nodes 5 and 7. This is because both
 325 substrings $\hat{S}_{i-1} \cdot \text{A}$ and $\hat{S}_{i-1} \cdot \text{B}$ occur in S . We can arbitrarily choose any of these nodes but,
 326 of course, the length of the substring we reconstruct may vary. Suppose we continue on node
 327 5. Then, we ask for the substring $\hat{S}_{i-1} \cdot \text{locus}(5)$. As the answer is positive, we continue with
 328 node 12 because $c_{12} = \text{B}$ and $\hat{S}_{i-1} \cdot \text{locus}(5) \cdot c_{12}$ occurs in S . We query for $\hat{S}_{i-1} \cdot \text{locus}(12)$. As
 329 the answer is negative, we binary search for the longest prefix P of $\text{locus}(12)$ such that $\hat{S}_{i-1} \cdot P$
 330 occurs in S . This way, we reconstruct the substring $X = \text{ABCAB}$ and we learn $\hat{S}_i = \hat{S}_{i-1} \cdot X$.

331 4. Conclusions and future work

332 We investigated the connection between the string reconstruction and compression prob-
 333 lems, establishing that they essentially coincide: the number of substring queries that need
 334 to be asked to an oracle in order to reconstruct a string S is proportional to the complexity
 335 of S .

336 We also showed that it is possible to efficiently reconstruct a string of length n over an
 337 alphabet of size σ using $O(\sigma g \log n) \subseteq O(\sigma \cdot z_{no} \log(n/z_{no}) \log n)$ queries in $O(n(\log \sigma + \log n))$
 338 time, where z_{no} is the number of phrases of the LZ77 factorization of S without overlaps and
 339 g is the size of the smallest grammar producing S . Immediate improvements over our work
 340 would be to replace z_{no} with the more powerful z (i.e., allowing overlaps), or to shave \log
 341 factors from the complexities of our algorithms. In particular, we know that the number of
 342 queries cannot be improved by more than a factor $O(\log \frac{n}{z_{no}} \log \sigma)$ in general.

343 In our setting, we aim at reconstructing the whole unknown string S . One can also
 344 consider the problem of reconstructing the set of all substrings of S of a given length k , see
 345 for example [4]. Notice that knowing all the substrings of S of length $r(S) + 2$ allows one

346 to uniquely determine S , where $r(S)$ is the repetition index of S , that is, the length of the
347 longest repeat of S [19, 20].

348 Another direction of investigation consists in introducing uncertainty into the model. For
349 example, allowing the oracle to answer the queries with a certain probability of returning a
350 wrong result — this could model strings with character ambiguities, e.g., DNA strings arising
351 from a sequencing — or allowing the oracle to return positive answers to queries within a
352 limited Hamming distance from substrings of the target string.

353 References

- 354 [1] T. Jiang, M. Li, DNA sequencing and string learning, *Mathematical systems theory*
355 29 (4) (1996) 387–405. doi:10.1007/BF01192694.
- 356 [2] D. Margaritis, S. Skiena, Reconstructing strings from substrings in rounds, in:
357 36th Annual Symposium on Foundations of Computer Science, Milwaukee, Wis-
358 consin, USA, 23-25 October 1995, IEEE Computer Society, 1995, pp. 613–620.
359 doi:10.1109/SFCS.1995.492591.
- 360 [3] S. Skiena, G. Sundaram, Reconstructing strings from substrings, *J. Comput. Biol.* 2 (2)
361 (1995) 333–353. doi:10.1089/cmb.1995.2.333.
- 362 [4] D. Tsur, Tight bounds for string reconstruction using substring queries, in: C. Chekuri,
363 K. Jansen, J. D. P. Rolim, L. Trevisan (Eds.), *Approximation, Randomization and Com-*
364 *binatorial Optimization, Algorithms and Techniques*, 8th International Workshop on
365 *Approximation Algorithms for Combinatorial Optimization Problems*, APPROX 2005
366 and 9th International Workshop on Randomization and Computation, RANDOM 2005,
367 Berkeley, CA, USA, August 22-24, 2005, Proceedings, Vol. 3624 of *Lecture Notes in*
368 *Computer Science*, Springer, 2005, pp. 448–459. doi:10.1007/11538462_38.
- 369 [5] J. Acharya, H. Das, O. Milenkovic, A. Orlitisky, S. Pan, String reconstruction
370 from substring compositions, *SIAM J. Discret. Math.* 29 (3) (2015) 1340–1371.
371 doi:10.1137/140962486.
- 372 [6] A. W. M. Dress, P. L. Erdős, Reconstructing words from subwords in linear time, *Annals*
373 *of Combinatorics* 8 (4) (2005) 457–462. doi:10.1007/s00026-004-0232-4.
- 374 [7] R. Afshar, A. Amir, M. T. Goodrich, P. Matias, Adaptive exact learning in a mixed-up
375 world: Dealing with periodicity, errors and jumbled-index queries in string reconstruc-
376 tion, in: *SPIRE 2020: Proceedings of the 27th International Symposium on String Pro-*
377 *cessing and Information Retrieval*, Vol. 12303 of *Lecture Notes in Computer Science*,
378 Springer, 2020, pp. 155–174. doi:10.1007/978-3-030-59212-7_12.
- 379 [8] K. Iwama, J. Teruyama, S. Tsuyama, Reconstructing strings from substrings: Optimal
380 randomized and average-case algorithms, *CoRR* abs/1808.00674 (2018).
381 URL <http://arxiv.org/abs/1808.00674>

- 382 [9] J. Ziv, A. Lempel, A universal algorithm for sequential data compression, *IEEE Trans-*
383 *actions on Information Theory* 23 (3) (1977) 337–343. doi:10.1109/TIT.1977.1055714.
- 384 [10] G. Navarro, Indexing highly repetitive string collections (2020). arXiv:2004.02781.
- 385 [11] M. Naor, String matching with preprocessing of text and pattern, in: *ICALP 1991:*
386 *Proceedings of the 18th International Colloquium on Automata, Languages, and Pro-*
387 *gramming*, Springer Berlin Heidelberg, Berlin, Heidelberg, 1991, pp. 739–750.
- 388 [12] M. A. Bender, M. Farach-Colton, B. C. Kuszmaul, Cache-oblivious string b-
389 trees, in: *PODS 2006: Proceedings of the Twenty-Fifth ACM SIGACT-SIGMOD-*
390 *SIGART Symposium on Principles of Database Systems*, 2006, pp. 233–242.
391 doi:10.1145/1142351.1142385.
- 392 [13] P. Ferragina, R. Venturini, Compressed cache-oblivious string B-tree, *ACM Transactions*
393 *on Algorithms (TALG)* 12 (4) (2016) 52:1–52:17. doi:10.1145/2903141.
- 394 [14] A. Amir, M. Farach, Y. Matias, Efficient randomized dictionary matching algorithms,
395 in: *Combinatorial Pattern Matching*, Springer Berlin Heidelberg, Berlin, Heidelberg,
396 1992, pp. 262–275.
- 397 [15] C. Jordan, Sur les assemblages de lignes, *Journal für die reine und angewandte Mathe-*
398 *matik* 70 (1869) 185–190.
- 399 [16] G. S. Brodal, R. Fagerberg, C. N. S. Pedersen, A. Östlin, The complexity of constructing
400 evolutionary trees using experiments, in: *ICALP 2001: Proceedings of the 28th Inter-*
401 *national Colloquium on Automata, Languages and Programming*, 2001, pp. 140–151.
402 doi:10.1007/3-540-48224-5_12.
- 403 [17] D. Della Giustina, N. Prezza, R. Venturini, A new linear-time algorithm for cen-
404 troid decomposition, in: *SPIRE 2019: Proceedings of the 26th International Sym-*
405 *posium on String Processing and Information Retrieval*, Springer, 2019, pp. 274–282.
406 doi:10.1007/978-3-030-32686-9_20.
- 407 [18] E. Ukkonen, On-line construction of suffix trees, *Algorithmica* 14 (3) (1995) 249–260.
408 doi:10.1007/BF01206331.
- 409 [19] A. Carpi, A. de Luca, Words and special factors, *Theoret. Comput. Sci.* 259 (1-2) (2001)
410 145–182. doi:10.1016/S0304-3975(99)00334-5.
- 411 [20] G. Fici, F. Mignosi, A. Restivo, M. Sciortino, Word assembly through minimal forbidden
412 words, *Theor. Comput. Sci.* 359 (1-3) (2006) 214–230. doi:10.1016/j.tcs.2006.03.006.