

# Regular Expression Matching on Compressed Text

Pierre Ganty

IMDEA Software Institute, Spain

pierre.ganty@imdea.org

Pedro Valero

IMDEA Software Institute, Spain & Universidad Politécnica de Madrid, Spain

pedro.valero@imdea.org

## Abstract

The state of the art approach to perform regular expression matching on compressed text is to feed the uncompressed text, as it is recovered by the decompressor, into the regular expression engine. We present **zearch**, a tool for searching *directly* on the compressed text, requiring linear time with respect to the size of the compressed data which might be exponentially smaller than its uncompressed version. The experiments show that **zearch**, in its current sequential implementation, outperforms some *highly optimized* implementations of the decompress and search approach even when decompressor and search engine each have a distinct computing unit.

**2012 ACM Subject Classification** F.4.2 Grammars and Other Rewriting Systems

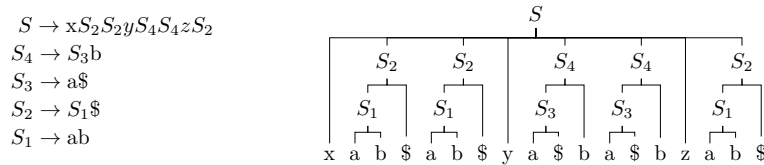
**Keywords and phrases** Straight Line Program, Grammar-based Compression, Regular Expression Matching

**Digital Object Identifier** 10.4230/LIPIcs.CVIT.2016.23

## 1 Introduction

The growing amount of information handled by modern systems demands efficient techniques both for compression, to reduce the storage cost, and for regular expression matching, to speed up processing. Despite the efforts invested in designing algorithms for searching with a regular expression *directly* in compressed text, there is no algorithm that nowadays outperforms feeding the output of the decompressor into the search engine.

Lossless compression of textual data is achieved by finding repetitions in the input text and replacing them by references. We focus on grammar-based compression schemes in which each tuple “reference  $\rightarrow$  repeated text” corresponds to a rule of a context-free grammar. The resulting grammar, produced as the output of the compression, generates a language consisting of a single word: the uncompressed text. Figure 1 depicts the output of a grammar-based compression algorithm.



**Figure 1** List of grammar rules (left) generating the string “xab\$ab\$ya\$ba\$zbab\$” (and no other) as evidenced by the parse tree (right).

Our approach is briefly described as follows: given a grammar-compressed text and a regular expression compute, for each variable of the grammar, the update on the state of the search that would result from processing the string generated from that variable (assuming the search is performed with an automata-based algorithm). The algorithm iterates through the grammar rules and composes, for each of them, the updates previously computed for the

variables on the right hand side. For example, when processing rule  $S_2 \rightarrow S_1 \$$  of Figure 1 our algorithm composes the update for  $S_1$  with the one for  $\$$ . The resulting information, i.e. the update that results from processing the string “ab\$”, will be reused every time the variable  $S_2$  appears in the right hand side of a rule.

We implement this approach in a tool called **zearch** that reports the number<sup>1</sup> of lines in the uncompressed text containing a substring that matches the expression. Our experiments show that **zearch** significantly outperforms previous tools for regular expression matching on compressed text and offers better results than some of the most efficient implementations of the decompress and search approach. The latter holds also when decompression and search are done in parallel, even though our implementation is purely sequential.

**Related work.** Results in regular expression matching on grammar-compressed text fall into two main groups: a) characterization of the problem’s complexity from a theoretical point of view [1, 2, 3] and b) development of algorithms and data structures to efficiently solve different instances of the problem (pattern matching [4, 5, 6], approximate pattern matching [7, 8], multi-pattern matching [9, 10], regular expression matching [11, 7]...).

Abboud et al. [3] proved that, under the Strong Exponential Time Hypothesis, there is no combinatorial algorithm deciding whether a grammar-compressed text contains a match of a deterministic FSA running on  $\mathcal{O}((p \cdot s)^{1-\varepsilon})$  time [Thm. 3.2] with  $\varepsilon > 0$  where  $p$  is the size of the compressed text and  $s$  is the number of states of the FSA. Similarly, for the case of non deterministic FSA they proved that, under the  $k$ -Clique Conjecture, there is no such algorithm running on  $\mathcal{O}((p \cdot s^3)^{1-\varepsilon})$  time [Thm. 4.2]. As we show in this paper, our algorithm is optimal for both scenarios since it exhibits  $\mathcal{O}(p \cdot s^3)$  and  $\mathcal{O}(p \cdot s)$  time complexities for the deterministic and non deterministic cases respectively.

On the other hand, algorithms for regular expression matching on grammar-compressed text are, traditionally, designed for a particular compression scheme [4, 11, 7]. The first algorithm to solve this problem is due to Navarro [11] and it is defined for LZ78/LZW compressed text. His algorithm reports all positions in the uncompressed text at which a substring that matches the expression ends and exhibits  $\mathcal{O}(2^s + s \cdot N + \text{occ} \cdot s \log s)$  worst case time complexity using  $\mathcal{O}(2^s + p \cdot s)$  space, where  $N$  is the size of the uncompressed text,  $\text{occ}$  is the number of occurrences, and  $p$  and  $s$  are defined as above. To the best of our knowledge this is the only algorithm for regular expression matching on compressed text that has been implemented and evaluated in practice. Bille et al. [7] improved the result of Navarro defining a data structure of size  $\mathcal{O}(p)$  to represent LZ78 compressed texts and an algorithm that, given a parameter  $\tau$ , finds all occurrences of a regular expression in a LZ78 compressed text in  $\mathcal{O}(p \cdot s(s + \tau) + \text{occ} \cdot s \cdot \log s)$  time using  $\mathcal{O}(p \cdot s^2/\tau + p \cdot s)$  space. Recall that  $\text{occ}$  might be as large as  $N$ .

We define an *occurrence* as a line of uncompressed text containing an a match of the expression and design an algorithm to report the number of occurrences in  $\mathcal{O}(p \cdot s^3 \cdot \log(N))$  time using  $\mathcal{O}(p \cdot s^2 \cdot \log(N))$  space. Hence, **zearch** reports the number of occurrences of a *fixed* regular expression in a compressed text in  $\mathcal{O}(p \cdot \log N)$  time while previous algorithms require  $\mathcal{O}(N)$  since  $\text{occ}$  might be as large as  $N$ . Note that when we are interested in printing the matches, the dependence with  $N$  in the running time cannot be avoided.

Finally, note that finding the smallest grammar –  $m^*$  – that generates a given string is an intractable problem and even approximating its size up to a small constant factor is NP-hard [12]. However there are several polynomial approximations such as LZ78 [13], LZW [14], RePair [15] or Sequitur [16] among others. While the algorithms of Navarro

<sup>1</sup> Optionally **zearch** also reports the matching lines.

and Bille et al. tie themselves to a particular approximation (LZ78/LZW) our approach is defined over straight line programs (grammars generating a single string), hence it applies to any grammar-based compression scheme, even for future ones yet to be defined. This generality of our algorithm has practical implications since, as shown by Hucke et al. [17], the LZ78 representation of a text of length  $N$  has size  $\Omega(|m^*| \times (N/\log N)^{2/3})$ , while its RePair representation has size  $\Omega(|m^*| \times \log N / \log \log N)$ . Therefore, even though the complexity of our algorithm and the one of Bille et al. depend on  $p$ , for the same uncompressed text we deal with much smaller values of  $p$ .

**Paper structure.** Section 2 introduces the notation used to prove<sup>2</sup>, in Section 3, the correctness of our algorithm; Section 4 describes the implementation of our tool while Section 5 discusses the experiments carried out and the obtained results; Section 5 analyzes the complexity of **zsearch** and Section 6 concludes with a brief discussion on some possible directions for future work.

## 2 Preliminaries

An *alphabet*  $\Sigma$  is a nonempty finite set of *symbols*. A *string*  $w$  is a finite sequence of symbols of  $\Sigma$  where the empty sequence is denoted  $\varepsilon$ . A *language* is a set of strings and the set of all strings over  $\Sigma$  is denoted  $\Sigma^*$ . We denote by  $|w|$  the *length* of  $w$  and denote it by  $\dagger$  when  $w$  is clear from the context. Further define  $(w)_i$  as the  $i$ -th symbol of  $w$  if  $1 \leq i \leq \dagger$  and  $\varepsilon$  otherwise. Similarly,  $(w)_{i,j}$  denotes the substring, also called *factor*, of  $w$  between the  $i$ -th and the  $j$ -th symbols, both included.

A *finite state automaton* (FSA or automaton for short) is a tuple  $A = (Q, \Sigma, I, F, \delta)$  where  $Q$  is a finite set of *states*;  $I \subseteq Q$  are the *initial states*;  $\Sigma$  is an alphabet;  $F \subseteq Q$  are the *final states*; and  $\delta \subseteq Q \times \Sigma \times Q$  are the *transitions*. Note that with this definition  $\varepsilon$ -transitions are not allowed. A *configuration* of an FSA is a pair  $(v, q)$  where  $v \in \Sigma^*$  represents the rest of the input and  $q \in Q$ , the current state. Define the *step* relation  $\rightarrow$  over configurations given by  $(v, q) \rightarrow (v', q')$  if  $v = a v'$  for some  $a \in \Sigma$  and  $(q, a, q') \in \delta$ . Given  $(v, q), (v', q')$  such that  $(v, q) \rightarrow^* (v', q')$  (the reflexive-transitive closure of  $\rightarrow$ ), there exists a witnessing sequence (possibly empty) of steps between  $(v, q)$  and  $(v', q')$  called a *run*. A run  $(v, q) \rightarrow^* (\varepsilon, q')$  is usually denoted  $q \xrightarrow{v} q'$ . The *language* of  $A$ , denoted  $\mathcal{L}(A)$ , is the set  $\{w \mid q_0 \xrightarrow{w} q, q_0 \in I, q \in F\}$ . A language  $L$  is said to be *regular*, if there exists an FSA  $A$  such that  $L = \mathcal{L}(A)$ . A string  $w$  *matches* against the automaton, and we say it is a *match*, if  $w \in \mathcal{L}(A)$ . For clarity, we assume through the paper that  $I$  and  $F$  are disjoint. Otherwise  $\varepsilon \in \mathcal{L}(A)$  therefore for every string there is a factor (in this case  $\varepsilon$ ) that matches the automaton.

A *context-free grammar* (or grammar for short) is a tuple  $G = (V, \Sigma, S, \mathcal{R})$  where  $V$  is a finite set of *variables* (or *non-terminals*) including the *axiom*  $S$ ;  $\Sigma$  is an alphabet (or set of *terminals*),  $\mathcal{R} \subseteq V \times (\Sigma \cup V)^*$  is a finite set of *rules*. Given a rule  $(X, \xi) \in \mathcal{R}$ , we often write it as  $X \rightarrow \xi$  for convenience. Define the *step* relation  $\Rightarrow$  as a binary relation over  $(V \cup \Sigma)^*$  given by  $\rho \Rightarrow \sigma$  if there exists a rule  $X \rightarrow \xi$  of  $G$  and a position  $i$  such that  $(\rho)_i = X$  and  $\sigma = (\rho)_{1,i-1} \xi (\rho)_{i+1,\dagger}$ . Given  $\rho, \sigma$  such that  $\rho \Rightarrow^* \sigma$  (the reflexo-transitive closure of  $\Rightarrow$ ), there exists a witnessing sequence (possibly empty) of steps between  $\rho$  and  $\sigma$ . Incidentally when  $\rho \in V$  and  $\sigma \in \Sigma^*$  such a step sequence is called a *derivation* (of variable  $\rho$ ). The *language* of  $G$ , denoted  $\mathcal{L}(G)$ , is the set  $\{w \in \Sigma^* \mid S \Rightarrow^* w\}$ . A language  $L$  is said to be *context-free*, or *CFL*, if there exists a grammar  $G$  such that  $L = \mathcal{L}(G)$ .

<sup>2</sup> Due to the lack of space some proofs are deferred to the Appendix.

A *Straight Line Program*  $P = (V, \Sigma, \mathcal{R})$ , hereafter SLP, is a grammar  $G = (V, \Sigma, S, \mathcal{R})$  such that  $\mathcal{R}$  is an indexed set of  $|V|$  rules (i.e.  $\mathcal{R} = (r_i)_{i \in 1, \dots, |V|}$ ) where each indexed rule has the form  $r_i = (X_i, \alpha_i \beta_i)$  with  $\alpha_i, \beta_i \in (\Sigma \cup \{X_j \in V \mid j < i\})$  and  $S = X_{|V|}$ . That is, variables are also indexed, there is exactly one rule per variable and their index coincide. Note that variables in the right hand side of  $r_i$  are left hand side of a rule with a lower index. We refer to  $X_{|V|} \rightarrow \alpha_{|V|} \beta_{|V|}$  as the *axiom rule*. We define the *size* of an SLP  $P$ , denoted  $|P|$ , as the number of rules,  $|\mathcal{R}|$ . It is easy to see that the language generated by an SLP consists of a single word  $w \in \Sigma^*$  and, by definition,  $|w| > 1$ . When  $\mathcal{L}(P) = \{w\}$  we abuse of notation and identify  $w$  with  $\mathcal{L}(P)$ .

An *Extended SLP*, hereafter ESLP, is a grammar  $(V \cup \{X_{|V|+1}\}, \Sigma, X_{|V|+1}, \mathcal{R} \cup \{r_{|V|+1}\})$  such that  $(V, \Sigma, \mathcal{R})$  is an SLP, and the *axiom rule*,  $r_{|V|+1} = (X_{|V|+1} \rightarrow \sigma)$ , has  $\sigma \in (\Sigma \cup V)^*$  with  $|\sigma| > 2$ . We define the size of an ESLP  $P$  as  $|P| = |\mathcal{R}| + |\sigma|$ .

### 3 Regular Expression Matching on Grammar-compressed Texts

Esparza et al. [18] defined the *saturation construction* as an algorithm to solve a number of decision problems involving automata and context-free grammars. In particular, given an automaton built from a regular expression and an SLP produced by a grammar-based compressor, the *saturation construction* can be used to decide whether or not the uncompressed text matches against the expression. We consider this algorithm as a starting point for searching with regular expressions in grammar-compressed texts.

By restricting the input grammar to SLPs, the *saturation construction* can be simplified since all rules have the same format and are listed in an orderly manner so that rule  $X \rightarrow \alpha\beta$  is always processed after  $\alpha$  and  $\beta$ . As a result we obtain Algorithm 0 which returns *true* iff the string generated by the SLP belongs to the language generated by the automaton, i.e. there is a run  $q_0 \xrightarrow{w} q_f$  with  $\mathcal{L}(P) = \{w\}$ ,  $q_0 \in I$  and  $q_f \in F$ .

---

#### Algorithm 0 Saturation construction for SLPs

---

**Input:** SLP  $P = (V, \Sigma, \mathcal{R})$ , FSA  $A = (Q, \Sigma, I, F, \delta)$

**Output:**  $\mathcal{L}(P) \stackrel{?}{\in} \mathcal{L}(A)$

```

1: function MAIN
2:   for each  $\ell = 1, 2, \dots, |V|$  do
3:     let  $(X_\ell \rightarrow \alpha_\ell \beta_\ell) \in \mathcal{R}$ 
4:     for each  $q_1, q' \in Q$  s.t.  $(q_1, \alpha_\ell, q') \in \delta$  do
5:       for each  $q_2 \in Q$  s.t.  $(q', \beta_\ell, q_2) \in \delta$  do
6:          $\delta := \delta \cup \{(q_1, X_\ell, q_2)\}$ 
7:   return  $((I \times \{X_{|V|}\} \times F) \cap \delta) \neq \emptyset$ 

```

---

However, standard search tools go beyond Algorithm 0 and check the existence of factors in  $w$  that match the regular expression. In our setup, this is equivalent to answering the question  $\mathcal{L}(P) \stackrel{?}{\in} (\Sigma^* \cdot \mathcal{L}(A) \cdot \Sigma^*)$  which can be done by adding  $\{(q, a, q) \mid q \in (I \cup F), a \in \Sigma\}$  to the set of transitions of the automaton before applying Algorithm 0. This will increase the number of operations carried out by the algorithm which will end up adding  $\{(q, X, q) \mid q \in (I \cup F), X \in V\}$  to  $\delta$ . To prevent the performance issues derived from handling these transitions, we consider them implicitly as shown by Algorithm 1.

---

#### Algorithm 1 Regex matching on SLP-compressed text

---

**Input:** SLP  $P = (V, \Sigma, \mathcal{R})$ , FSA  $A = (Q, \Sigma, I, F, \delta)$

**Output:**  $\mathcal{L}(P) \stackrel{?}{\in} (\Sigma^* \cdot \mathcal{L}(A) \cdot \Sigma^*)$

```

1: function MAIN
2:   for each  $\ell = 1, 2, \dots, |V|$  do
3:     let  $(X_\ell \rightarrow \alpha_\ell \beta_\ell) \in \mathcal{R}$ 
4:     for each  $q_1, q' \in Q$  s.t.  $(q_1, \alpha_\ell, q') \in \delta$  or  $q_1 = q' \in I$  do
5:       for each  $q_2 \in Q$  s.t.  $(q', \beta_\ell, q_2) \in \delta$  or  $q' = q_2 \in F$  do
6:          $\delta := \delta \cup \{(q_1, X_\ell, q_2)\}$ 
7:   return  $((I \times \{X_{|V|}\} \times F) \cap \delta) \neq \emptyset$ 

```

---

► **Theorem 1.** Let  $P = (V, \Sigma, \mathcal{R})$  be an SLP and  $A = (Q, \Sigma, I, F, \delta)$  an FSA. Let  $X \in V$ ,  $w \in \Sigma^*$  with  $X \Rightarrow^* w$ . After  $\ell$  iterations of the outermost “for” loop of Algorithm 1, if  $X \in \{X_1, \dots, X_\ell\} \subseteq V$  then  $((I \times \{X\} \times F) \cap \delta) \neq \emptyset \iff w \in \Sigma^* \cdot \mathcal{L}(A) \cdot \Sigma^*$ .

153 ► **Corollary 2.** Let  $P = (V, \Sigma, \mathcal{R})$  be an SLP and  $A = (Q, \Sigma, I, F, \delta)$  an FSA. Algorithm 1  
 154 returns true iff  $\mathcal{L}(P) \in (\Sigma^* \cdot \mathcal{L}(A) \cdot \Sigma^*)$ .

### 155 3.1 Counting Matching Lines

156 State of the art tools for regular expression search are equipped with a number of features<sup>3</sup>  
 157 to perform different operations beyond deciding the existence of a match in the text. Among  
 158 the most relevant of these features we find *counting*. Tools like **grep**<sup>4</sup>, **rg**<sup>5</sup>, or **ack**<sup>6</sup> report  
 159 the number of lines containing a match, ignoring matches across lines. Next we modify  
 160 Algorithm 1 to count matching lines, denoting the new-line delimiter by symbol  $\hookrightarrow$  and  
 161  $\widehat{\Sigma} = \Sigma \setminus \{\hookrightarrow\}$ . Given  $w \in \Sigma^*$ , a *line* is a maximal factor of  $w$  each symbol of which belongs  
 162 to  $\widehat{\Sigma}$ . A *closed line* is a line which is not a prefix nor a suffix of  $w$ .

163 ► **Definition 3** (Matching lines). Let  $P=(V, \Sigma, \mathcal{R})$  be an SLP with  $\tau \in (V \cup \Sigma)$ ,  $w \in \Sigma^+$  and  
 164  $\tau \Rightarrow^* w$ . Let  $A=(Q, \widehat{\Sigma}, I, F, \delta)$  be an FSA. The *matching lines generated by  $\tau$*  is the set

$$165 \quad \mathcal{ML}_\tau = \left\{ (i, j) \mid (w)_{i,j} \in \widehat{\Sigma}^* \cdot \mathcal{L}(A) \cdot \widehat{\Sigma}^*, (i=1 \vee (w)_{i-1}=\hookrightarrow), (j=|w| \vee (w)_{j+1}=\hookrightarrow) \right\}.$$

166 ► **Example 4.** Let  $X$  be a variable with  $X \Rightarrow^* ba\hookrightarrow ab\hookrightarrow aba$  and let  $A$  be an automaton with  
 167  $\mathcal{L}(A) = \{ab, ba\}$ . Then  $\mathcal{ML}_X = \{(1, 2), (4, 5), (7, 9)\}$   $\diamond$

168 In order to count the matching lines generated by an SLP with axiom  $X_{|V|}$ , it suffices to  
 169 compute  $|\mathcal{ML}_{X_{|V|}}|$ . Lemmas 5 and 6 show that given a rule  $X \rightarrow \alpha\beta$  the set  $\mathcal{ML}_X$  can be  
 170 computed using  $\mathcal{ML}_\alpha$ ,  $\mathcal{ML}_\beta$  and information about the transitions added by Algorithm 1.

171 ► **Lemma 5.** Let  $P = (V, \Sigma, \mathcal{R})$  be an SLP and  $A = (Q, \widehat{\Sigma}, I, F, \delta)$  an FSA. Let  $u, v \in \Sigma^+$   
 172 such that  $X \Rightarrow \alpha\beta \Rightarrow^* u \cdot \beta \Rightarrow^* u \cdot v$ .

- 173 1.  $\forall 1 \leq i \leq j < |u|, ((i, j) \in \mathcal{ML}_\alpha \iff (i, j) \in \mathcal{ML}_X)$ .
- 174 2.  $\forall 1 < i \leq j \leq |v|, ((i, j) \in \mathcal{ML}_\beta \iff (i+|u|, j+|u|) \in \mathcal{ML}_X)$ .

175 ► **Lemma 6.** Let  $P = (V, \Sigma, \mathcal{R})$  be an SLP and  $A = (Q, \widehat{\Sigma}, I, F, \delta)$  an FSA. Let  $u, v \in \Sigma^+$  such  
 176 that  $X \Rightarrow \alpha\beta \Rightarrow^* u \cdot \beta \Rightarrow^* u \cdot v$ . Then there exists  $(\ell_X, r_X) \in \mathcal{ML}_X$  such that  $\ell_X \leq |u| < r_X$   
 177 if and only if one of the following holds:

- 178 1.  $(\ell_X, |u|) \in \mathcal{ML}_\alpha$ .
  - 179 2.  $(1, r_X - |u|) \in \mathcal{ML}_\beta$ .
  - 180 3. Algorithm 1 adds  $(q_0, \alpha, q')$ ,  $(q', \beta, q_f)$  to  $\delta$  with  $q_0 \in I$ ,  $q' \notin (I \cup F)$  and  $q_f \in F$ .
- 181 Furthermore, when the pair  $(\ell_X, r_X)$  exists it is unique and

$$182 \quad \ell_X = \max\{0, i \mid (u)_i = \hookrightarrow\} + 1 \text{ and } r_X = |u| + \min\{|v| + 1, i \mid (v)_i = \hookrightarrow\} - 1.$$

183 ► **Example 7.** Let  $A$  be an automaton with  $\mathcal{L}(A) = \{ab, ba\}$ . Consider the grammar rule  
 184  $X \rightarrow \alpha\beta$  with  $\alpha \Rightarrow^* ba\hookrightarrow a$  and  $\beta \Rightarrow^* b\hookrightarrow aba$ . Then  $X \Rightarrow^* ba\hookrightarrow ab\hookrightarrow aba$ . By Definition 3 we  
 185 have  $\mathcal{ML}_\alpha = \{(1, 2)\}$  and  $\mathcal{ML}_\beta = \{(3, 5)\}$ .

186 Due to properties 1 and 2 of Lemma 5 we have  $(1, 2), (7, 9) \in \mathcal{ML}_X$  and due to property  
 187 3 of Lemma 6 we have  $(4, 5) \in \mathcal{ML}_X$ . Therefore  $\mathcal{ML}_X = \{(1, 2), (4, 5), (7, 9)\}$ .  $\diamond$

188 Observe that, for an uncompressed text of size  $N$ ,  $\mathcal{ML}_\tau$  has up to  $N$  elements yet its  
 189 size can be represented in space  $\log N$ . Next we define the *counting information* of a symbol  
 190  $\tau$ , denoted  $\mathcal{C}_\tau$ , of size logarithmic in  $\mathcal{ML}_\tau$ . For a variable  $X$  generating the string  $w$ ,  $\mathcal{C}_X$

<sup>3</sup> <https://beyondgrep.com/feature-comparison/>

<sup>4</sup> <https://www.gnu.org/software/grep>

<sup>5</sup> <https://github.com/BurntSushi/ripgrep>

<sup>6</sup> <https://github.com/beyondgrep/ack2>

contains: the existence of end-of-line delimiters in  $w$  ( $\text{NL}(X)$ ); the existence of a prefix (resp. suffix) of  $w$  containing a match ( $\text{L}(X)$ , resp.  $\text{R}(X)$ ) and the number of closed lines in  $w$  containing a match ( $\text{M}(X)$ ). Below Lemma 10 shows that, given a rule  $X \rightarrow \alpha\beta$ ,  $\mathcal{C}_X$  can be computed using  $\mathcal{C}_\alpha$ ,  $\mathcal{C}_\beta$  and the information about the transitions added by Algorithm 1.

► **Definition 8** (Counting information). Let  $P = (V, \Sigma, \mathcal{R})$  be an SLP and  $A = (Q, \widehat{\Sigma}, I, F, \delta)$  an FSA. Let  $\tau \in (V \cup \Sigma)$  and  $w \in \Sigma^+$  with  $\tau \Rightarrow^* w$  and

$$\begin{aligned} \text{L}(\tau) &:= \exists i \text{ s.t. } (0, i) \in \mathcal{ML}_\tau & \text{R}(\tau) &:= \exists j \text{ s.t. } (j, \dagger) \in \mathcal{ML}_\tau \\ \text{NL}(\tau) &:= \exists k \text{ s.t. } (w)_k = \lhd & \text{M}(\tau) &:= |\{(i, j) \in \mathcal{ML}_\tau \mid 1 < i \leq j < |w|\}| \end{aligned}$$

We define the *counting information* of  $\tau$  as  $\mathcal{C}_\tau = \langle \text{NL}(\tau), \text{L}(\tau), \text{R}(\tau), \text{M}(\tau) \rangle$ .

► **Remark 9.** Since  $\lhd \notin \widehat{\Sigma}$ , pairs in  $\mathcal{ML}_\tau$  do not overlap. Hence

$$|\mathcal{ML}_\tau| = \text{M}(\tau) + \begin{cases} 1 & \text{if } \text{L}(\tau) \vee \text{R}(\tau) \\ 0 & \text{otherwise} \end{cases} + \begin{cases} 1 & \text{if } \text{NL}(\tau) \wedge \text{L}(\tau) \wedge \text{R}(\tau) \\ 0 & \text{otherwise} \end{cases}$$

As mentioned before, to count the number of matching lines generated by an SLP with axiom  $X_{|V|}$ , it suffices to compute  $|\mathcal{ML}_{X_{|V|}}|$ . By Remark 9, this amounts to compute  $\mathcal{C}_{X_{|V|}}$ .

► **Lemma 10.** Let  $P = (V, \Sigma, \mathcal{R})$  be an SLP and  $A = (Q, \widehat{\Sigma}, I, F, \delta)$  an FSA. For each rule  $(X \rightarrow \alpha\beta) \in \mathcal{R}$  let  $m = \text{true}$  iff Algorithm 1 adds  $(q_0, \alpha, q')$ ,  $(q', \beta, q_f)$  to  $\delta$  with  $q_0 \in I$ ,  $q' \notin (I \cup F)$  and  $q_f \in F$ . The following equalities hold:

1.  $\text{NL}(X) = \text{NL}(\alpha) \vee \text{NL}(\beta)$
2.  $\text{L}(X) = \begin{cases} \text{L}(\alpha) \vee \text{L}(\beta) \vee m & \text{if } \neg \text{NL}(\alpha) \\ \text{L}(\alpha) & \text{otherwise} \end{cases}; \quad \text{R}(X) = \begin{cases} \text{R}(\alpha) \vee \text{R}(\beta) \vee m & \text{if } \neg \text{NL}(\beta) \\ \text{R}(\beta) & \text{otherwise} \end{cases}$
3.  $\text{M}(X) = \text{M}(\alpha) + \text{M}(\beta) + \begin{cases} 1 & \text{if } \text{NL}(\alpha) \wedge \text{NL}(\beta) \wedge (\text{R}(\alpha) \vee \text{L}(\beta) \vee m) \\ 0 & \text{otherwise} \end{cases}$

► **Example 11.** Let  $A$  be an automaton with  $\mathcal{L}(A) = \{ab, ba\}$ . Consider the grammar rule  $X \rightarrow \alpha\beta$  with  $\alpha \Rightarrow^* ba\lhd a$  and  $\beta \Rightarrow^* b\lhd aba$ . Then  $X \Rightarrow^* ba\lhd ab\lhd aba$ . By Definition 3 we have  $\mathcal{C}_\alpha = \langle 1, 1, 0, 0 \rangle$  and  $\mathcal{C}_\beta = \langle 1, 0, 1, 0 \rangle$ .

Following Lemma 10 we find that  $\mathcal{C}_X = \langle 1, 1, 1, 1 \rangle$ , hence  $|\mathcal{ML}_X| = 1 + 1 + 1 = 3$ , by Remark 9, which coincides with the length of the set  $\mathcal{ML}_X$  computed in Example 7. ◊

Algorithm 2 extends Algorithm 1 using the counting information to report the number of lines of uncompressed text containing a match of the regular expression. We use the conditional operator (**cond** ? **op1** : **op2**) that returns **op1** if **cond** is *true* and **op2** otherwise. For each symbol  $\tau$ ,  $\langle \tau_{\text{NL}}, \tau_{\text{L}}, \tau_{\text{R}}, \tau_{\text{M}} \rangle$  denotes the information computed and stored by Algorithm 2 for  $\tau$ . As Theorem 12 shows, this equals to the *counting information* of each symbol  $\tau$ .

► **Theorem 12.** Let  $P = (V, \Sigma, \mathcal{R})$  be an SLP and  $A = (Q, \widehat{\Sigma}, q_0, F, \delta)$  an FSA. Algorithm 2 returns  $|\mathcal{ML}_{X_{|V|}}|$ .

**Proof.** Note that Definitions 3 and 8 applied to terminals  $a \in \Sigma$  result in  $\mathcal{ML}_a$  either empty or containing a single pair  $(1, 1)$  and:

$$\text{L}(a) = \text{R}(a) = \begin{cases} 1 & \text{if } (q_0, a, q_f) \in \delta, q_f \in F \\ 0 & \text{otherwise} \end{cases} \quad \text{NL}(a) = \begin{cases} 1 & \text{if } a = \lhd \\ 0 & \text{otherwise} \end{cases}.$$

It is straightforward to observe that a) after running function `INIT_AUTOMATON` of Algorithm 2,  $a_{\text{NL}} = \text{NL}(a)$ ,  $a_{\text{L}} = a_{\text{R}} = \text{L}(a) = \text{R}(a)$  and  $a_{\text{M}} = \text{M}(a)$  for each  $a \in \Sigma$  and b) if  $\alpha_{\text{NL}} = \text{NL}(\alpha)$ ,  $\alpha_{\text{L}} = \text{L}(\alpha)$ ,  $\alpha_{\text{R}} = \text{R}(\alpha)$  and  $\alpha_{\text{M}} = \text{M}(\alpha)$  (and the same for  $\beta$ ) then after running function `COUNT` of Algorithm 2 equalities from Lemma 10 hold when replacing  $\langle \text{NL}(\tau), \text{L}(\tau), \text{R}(\tau), \text{M}(\tau) \rangle$  by  $\langle \tau_{\text{NL}}, \tau_{\text{L}}, \tau_{\text{R}}, \tau_{\text{M}} \rangle$ . Furthermore, Theorem 1 and Lemmas 5, 6 and 10

**Algorithm 2** Extension of Algorithm 1 to report the number of matching lines

---

**Input:** An SLP  $P = (V, \Sigma, \mathcal{R})$  and an FSA  $A = (Q, \widehat{\Sigma}, I, F, \delta)$   
**Output:** Number of lines of  $\mathcal{L}(P)$  containing a factor matching  $A$ .

```

1: procedure COUNT( $X, \alpha, \beta, m$ )
2:    $X_{\text{NL}} := \alpha_{\text{NL}} \vee \beta_{\text{NL}}$ 
3:    $X_{\text{L}} := (\neg \alpha_{\text{NL}} ? \alpha_{\text{L}} \vee \beta_{\text{L}} \vee m : \alpha_{\text{L}}); \quad X_{\text{R}} := (\neg \beta_{\text{NL}} ? \alpha_{\text{R}} \vee \beta_{\text{R}} \vee m : \beta_{\text{R}})$ 
4:    $X_{\text{M}} := \alpha_{\text{M}} + \beta_{\text{M}} + (\alpha_{\text{NL}} \wedge \beta_{\text{NL}} \wedge (\alpha_{\text{R}} \vee \beta_{\text{L}} \vee m) ? 1 : 0)$ 

5: procedure INIT_AUTOMATON( )
6:   for each  $a \in \Sigma$  do
7:      $a_{\text{NL}} := (a = \lhd)$ 
8:      $a_{\text{L}} := ((q_0, \alpha, q_f) \in \delta, q_0 \in I, q_f \in F); \quad a_{\text{R}} := a_{\text{L}}$ 
9:      $a_{\text{M}} := 0$ 

10: function MAIN
11:   INIT_AUTOMATON( )
12:   for each  $\ell = 1, 2, \dots, |V|$  do
13:     let  $(X_\ell \rightarrow \alpha_\ell \beta_\ell) \in \mathcal{R}$ 
14:      $\text{new\_match} := \text{false}$ 
15:     for each  $q_1, q' \in Q$  s.t.  $(q_1, \alpha_\ell, q') \in \delta$  or  $q_1 = q' \in I$  do
16:       for each  $q_2 \in Q$  s.t.  $(q', \beta_\ell, q_2) \in \delta$  or  $q' = q_2 \in F$  do
17:          $\delta := \delta \cup \{(q_1, X_\ell, q_2)\}$ 
18:          $\text{new\_match} := \text{new\_match} \vee (q_1 \in I \wedge q' \notin (I \cup F) \wedge q_2 \in F)$ 
19:     COUNT( $X_\ell, \alpha_\ell, \beta_\ell, \text{new\_match}$ )
20:   return  $(X_{|V|})_{\text{M}} + ((X_{|V|})_{\text{NL}} ? (X_{|V|})_{\text{L}} + (X_{|V|})_{\text{R}} : (X_{|V|})_{\text{L}})$ 

```

---

also apply to Algorithm 2 since it performs *exactly* the same operations on  $\delta$  as Algorithm 1.  
 Therefore Algorithm 2 propagates the counting information from the terminals to the axiom,  
 Finally, by Remark 9 the value returned by the algorithm is  $|\mathcal{ML}_{X_{|V|}}|$ . ◀

## 4 Implementation

In this section we describe **zearch**, an implementation based on Algorithm 2 to perform regular expression matching on grammar-compressed text. This tool works on texts compressed with **repair**<sup>7</sup>, which implements the Recursive Pairing algorithm defined by Larsson et al. [15]. The choice of a particular compressor simplifies the implementation task since **zearch** has to read the grammar rules from the compressed file. However **zearch** can handle any grammar-based scheme by providing a way to recover the SLP from the input file. We use the automata library **libfa**<sup>8</sup> to translate the input regular expression into an  $\varepsilon$ -free FSA. Internally it uses Thompson's algorithm [19] with on-the-fly  $\varepsilon$ -removal.

### 4.1 Data structures

We index the transitions of the automaton by their label, i.e. we store a pointer to the list of labeled transitions, together with the counting information of each symbol, in an array where the  $i$ -th position contains the information related to the  $i$ -th symbol of the grammar. This representation of the automaton allows the loops in lines 15 and 16 to iterate directly over the transitions labeled by  $\alpha$  and  $\beta$  respectively.

<sup>7</sup> <https://storage.googleapis.com/google-code-archive-downloads/v2/code.google.com/re-pair/repair110811.tar.gz>

<sup>8</sup> <http://augeas.net/libfa/index.html>



This data structure causes the loop in line 16 of Algorithm 2 to iterate through all transitions labeled with  $\beta_\ell$ , checking whether they leave from  $q'$  or not. We improve this situation by using an auxiliary matrix with  $s^2$  elements of size  $\log s$  bits. Before entering the loop in line 15, **zearch** stores the elements of the set  $\beta_\ell(i) = \{q_j \mid (q_i, \beta_\ell, q_j) \in \delta\} \subseteq Q$  at positions  $(i, 1) \dots (i, |\beta_\ell(i)|)$  of the matrix. Whenever  $|\beta_\ell(i)| < s$ , it also stores  $-1$  at position  $(i, \beta_\ell(i)+1)$  to mark the end of the set. This allows the loop in line 16 of Algorithm 2 to iterate directly over the transitions labeled with  $\beta_\ell$  leaving from state  $q'$ .

To have constant time addition to the list of transitions labeled by a variable while avoiding repetitions, **zearch** requires an auxiliary matrix with  $s^2$  elements each of which of size  $\log p$  bits. When transition  $(q_i, X_\ell, q_j)$  is added to  $\delta$ , **zearch** first checks if position  $(i, j)$  of the matrix contains the value  $\ell$ . If not, then it adds the pair  $(q_i, q_j)$  to the list of transitions labeled by  $X_\ell$  and stores  $\ell$  in the matrix at position  $(i, j)$ . Note that the matrix can be reused for all rules since they are processed one at a time in an orderly manner.

Finally, the *counting information* of each variable  $X$  can be stored using one bit for each  $X_{\text{NL}}$ ,  $X_{\text{L}}$  and  $X_{\text{R}}$  and  $\log N$  bits for  $X_{\text{M}}$ , where  $N$  is the size of the uncompressed text. Hence, both function COUNT and INIT\_AUTOMATON are executed in  $\mathcal{O}(\log N)$  time.

These data structures allow Algorithm 2 to exhibit  $\mathcal{O}(p \cdot s^3 \cdot \log N)$  time complexity using  $\mathcal{O}(p \cdot s^2 \cdot \log N)$  space for the worst case scenario and  $\mathcal{O}(p \cdot \log N + s)$  time complexity using  $\mathcal{O}(p \cdot \log N + s)$  space for the best one<sup>9</sup>. Further details on the complexity of **zearch** are given in Section 5.

**Previously proposed data structures.** The data structure suggested by Esparza et al. [18] for the *saturation construction*, in which Algorithm 0 is inspired, consists of a bit array where each position represents a possible transition  $(q_i, X_\ell, q_j)$  and the value thereof indicates whether the transition is present in the automaton. This data structure offers constant time membership test and addition to the set of transitions using  $\mathcal{O}(ps^2)$  space which makes it a good candidate for implementing Algorithm 2. However, this representation of the automaton forces the loop in line 15 to perform  $s^2$  iterations since it has to consider all pairs of states and check whether they are connected by a transition labeled with  $\alpha$ . Similarly, the loop in line 16 requires  $s$  iterations regardless of the number of transitions labeled with  $\beta$ . Furthermore, in order to perform the counting, this data structure should be extended with an array of size  $\mathcal{O}(p \cdot \log N)$  for keeping the counting information, which should also be updated each time a rule is processed. As a result, with this data structure, Algorithm 2 operates in  $\Theta(p \cdot s^3 \cdot \log N)$  time and  $\Theta(p \cdot s^2 \cdot \log N)$  space.

## 4.2 Processing the axiom rule.

Grammar-based compression algorithms typically produce extended straight line programs. This happens because the axiom rule is built with the remains of the input text after extracting all repeated factors and replacing them by variables of the grammar. Once the algorithm stops there is no repeated factor in the remaining text, otherwise it would result in a new rule. The remaining text is then folded into an axiom rule which typically contains more symbols than rules are in the grammar so the way in which the axiom is handled has a great impact in the performance.

An ESLP can be transformed into an SLP by rewriting the axiom rule  $X_{|V|} \rightarrow \sigma$  as  $\{S_1 \rightarrow (\sigma)_1 (\sigma)_2\} \cup \{S_i \rightarrow S_{i-1}(\sigma)_{i+1} \mid i = 2 \dots |\sigma|-2\} \cup \{X_{|V|} \rightarrow S_{|\sigma|-2}(\sigma)_\dagger\}$  so Algorithm 2 could process the axiom rule by translating it into an SLP. It is worth pointing that the transitions labeled  $S_i$  can be discarded after processing the grammar rule with  $S_{i+1}$  on the

<sup>9</sup> The best case scenario corresponds to the loops in lines 15 and 16 iterating over empty sets.



left hand side. Hence we can reuse the array entry for  $S_i$  when processing  $S_{i+1}$ . Similarly, since for any  $S_i$  there is exactly one derivation  $X_{|V|} \Rightarrow^* \gamma S_i \rho$  and  $\gamma = \varepsilon$ , it suffices to store the transitions labeled by  $S_i$  leaving from the initial state. As a result we obtain Algorithm 3 (see Appendix B) which exhibits  $\mathcal{O}(\log N ((p - |\sigma|) \cdot s^3 + |\sigma| \cdot s^2))$  worst case time complexity using  $\mathcal{O}((p - |\sigma|) \cdot s^2 \cdot \log N)$  space where  $p$  is the size of the ESLP and  $X_{|V|} \rightarrow \sigma$  its axiom.

### 4.3 Printing matching lines

The counting information computed by Algorithm 2 can be used to *partially* decompress the input file, *fully decompressing* the matching lines and *skipping* symbols generating no matching lines. This partial decompression of the input file requires **zearch** to store the grammar rules, indexed by the symbol on the left hand side, which can be done in  $\mathcal{O}(p)$  space. Therefore, **zearch** has the same space complexity regardless of whether the matching lines are being reported or simply counted.

On the other hand, the time complexity does increase when **zearch** is asked to report the matching lines. In particular, the partial decompression described above requires  $\mathcal{O}(v \cdot d)$  extra time, where  $v$  is the size of the output generated and  $d$  is the depth of the grammar, i.e. the length of the longest sequence  $X_{i_1} X_{i_2} \dots X_{i_d}$  of variables of the SLP such that  $(X_{i_j} \rightarrow \alpha X_{i_{j+1}} \beta) \in \mathcal{R}$  with  $\alpha, \beta \in (\Sigma \cup V \cup \{\varepsilon\})$ . Note that this is the running time required to recover the uncompressed matching lines from the compressed file.

## 5 Empirical Evaluation

We measure the performance of our implementation by comparing it against **LZgrep** [4] and **GNgrep** [11], described in the *Related Work* section. Both tools operate on LZW-compressed text and, while **LZgrep** is restricted to pattern matching, **GNgrep** performs regular expression matching. These are, to the best of our knowledge, the only two existing tools for regular expression matching on compressed text.

As shown by Table 2, **zearch** significantly outperforms these two tools even by an order of magnitude in some cases. Therefore, we also compared its performance against the state of the art technique for regular expression matching on compressed text: feeding the output of the decompressor into the search engine. We implement this approach by decompressing the file with **zstd**<sup>10</sup> or **lz4**<sup>11</sup> and searching with **grep** or **rg**. These tools are top of the class for lossless compression<sup>12</sup> and regular expression matching<sup>13</sup>, respectively. We use **zstd** and **lz4** with the highest compression level since it has negligible impact in the time required for decompression. We also consider **gzip** (with highest compression level) as a reference.

For the experiments we consider versions **grep** v3.1, **rg** v0.7.1, **lz4** v1.8.2, **zstd** v1.3.2, and **gzip** v1.9 running in an Intel Xeon E5640 CPU 2.67 GHz with 20 GB RAM.

**Design of the experiments.** The benchmark designed to assess the performance of **zearch** combines ideas from those for regular expression engines and compression tools. As such, we consider different sizes and formats for the input files since both factors have an impact in the result of the compression. In particular our benchmark includes files containing English *subtitles* extracted from the OpenSubtitles2016 dataset [20], a sample of the public timeline recorded by GitHub Archive<sup>14</sup> in *JSON* format, a sample of The Google Books Ngram

<sup>10</sup><https://github.com/facebook/zstd>

<sup>11</sup><https://github.com/lz4/lz4>

<sup>12</sup><https://quixdb.github.io/squash-benchmark/>

<sup>13</sup><https://rust-leipzig.github.io/regex/2017/03/28/comparison-of-regex-engines/>

<sup>14</sup><https://www.gharchive.org/>

Viewer dataset<sup>15</sup> considering n-grams from 5 to 20 characters long in *CSV* format and a *Log* file containing two months worth of all HTTP requests to the NASA Kennedy Space Center WWW server in Florida<sup>16</sup>. Table 1 shows how each compressor behaves on these files<sup>17</sup>.

	File	Compressed size					Compression time					Decompression time				
		LZW	gzip	repair	zstd	lz4	LZW	gzip	repair	zstd	lz4	LZW	gzip	repair	zstd	lz4
Uncompressed size 100 MB	Logs	<b>0.19</b>	0.1	<b>0.08</b>	<b>0.07</b>	0.12	<b>0.04</b>	0.04	0.19	<b>0.56</b>	<b>0.03</b>	<b>0.02</b>	0.01	0.01	<b>0.004</b>	<b>0.003</b>
	Subtitles	<b>0.36</b>	0.33	<b>0.13</b>	<b>0.11</b>	0.15	<b>0.04</b>	0.1	<b>0.25</b>	0.21	<b>0.03</b>	<b>0.02</b>	0.01	0.01	<b>0.004</b>	<b>0.003</b>
	JSON	<b>0.27</b>	0.14	<b>0.12</b>	<b>0.1</b>	0.17	0.05	<b>0.03</b>	0.22	<b>0.4</b>	<b>0.02</b>	<b>0.02</b>	0.01	0.01	<b>0.004</b>	<b>0.004</b>
	CSV	0.37	0.38	<b>0.33</b>	<b>0.3</b>	<b>0.41</b>	<b>0.04</b>	0.21	0.26	<b>0.43</b>	<b>0.11</b>	<b>0.02</b>	0.01	0.02	<b>0.01</b>	<b>0.004</b>
	Logs	<b>19</b>	11	<b>6.8</b>	<b>6.3</b>	12	<b>3.4</b>	3.7	25.9	<b>116.4</b>	<b>2.8</b>	<b>1.4</b>	0.5	1.0	<b>0.18</b>	<b>0.08</b>
	Subtitles	<b>39</b>	31	<b>14</b>	<b>11</b>	23	<b>4.0</b>	9.5	31.6	<b>62.0</b>	<b>4.8</b>	<b>1.6</b>	0.83	1.5	<b>0.2</b>	<b>0.08</b>
	JSON	<b>26</b>	14	<b>8.8</b>	<b>6.8</b>	16	<b>3.7</b>	<b>3.2</b>	25.7	<b>59.5</b>	<b>2.0</b>	<b>1.4</b>	0.58	1.1	<b>0.17</b>	<b>0.08</b>
	CSV	37	37	<b>27</b>	<b>24</b>	<b>40</b>	<b>3.9</b>	21.6	31.2	<b>111.7</b>	<b>10.6</b>	<b>1.6</b>	0.93	1.6	<b>0.35</b>	<b>0.1</b>

Table 1 Sizes (in MB) of the compressed files and (de)compression times (in seconds). Values measured with maximum compression levels enabled. (Blue = best tool; bold black text = second best; red = worst one).

We first run each experiment 3 times as warm up so that the files are loaded in memory, reducing the impact of the I/O bound. Then we measure the running time 30 times and compute the *confidence interval* (with 95% confidence) for the average running time required to count the number of lines matching a regular expression in a certain file using a certain tool. Table 2 summarizes the obtained results when considering, for all files, the regular expressions<sup>18</sup>: “qwerty”, “Hello”, “.”, “I . \* you”, “ [a-z]{4} ” “ [a-z]\*[a-z]{6} ”, “ [0-9]7[0-9]4[0-9]9[0-9]0[0-9] ”, “([a-z]{5} +){4}”. A more extensive evaluation considering a bigger range of file sizes and regular expressions is available online<sup>19</sup>.

Tool	Logs	Subtitles	JSON	CSV	Logs	Subtitles	JSON	CSV
grep	1.02	0.91	0.67	0.51	1.57	0.98	0.70	0.49
rg	1.59	1.57	0.82	1.00	2.38	1.56	0.45	1.16
zsearch	<b>4.08</b>	<b>6.41</b>	<b>5.36</b>	9.47	<b>118.77</b>	<b>350.97</b>	<b>180.11</b>	509.54
GNgrep	13.55	9.70	10.88	6.54	6.83	4.41	5.73	2.69
LZgrep	3.30	1.86	2.54	1.24	9.12	2.41	5.93	1.67
gzip+grep	2.60 (1.90)	2.48 (1.76)	1.98 (1.59)	1.62 (1.26)	6.09 (4.43)	3.41 (2.41)	4.02 (3.28)	2.32 (1.82)
lz4+grep	1.66 (1.34)	1.29 (1.10)	1.17 (0.91)	<b>0.82</b> (0.68)	2.55 (2.19)	1.31 (1.18)	1.29 (1.11)	<b>0.74</b> (0.66)
lz4+rg	2.30 (1.77)	2.01 (1.68)	1.48 (1.03)	1.37 (1.10)	3.80 (3.34)	2.05 (1.88)	1.51 (1.23)	1.50 (1.40)
zstd+grep	1.80 (1.33)	1.41 (1.08)	1.30 (1.01)	1.02 (0.78)	3.41 (2.36)	1.66 (1.20)	1.84 (1.49)	1.24 (0.97)
zstd+rg	2.45 (1.94)	2.12 (1.75)	1.61 (1.24)	1.57 (1.23)	4.62 (3.45)	2.40 (1.88)	1.97 (1.57)	2.00 (1.55)

Table 2 Running times (in milliseconds) required to search for a regular expression in a compressed file of uncompressed size 1 MB (left) and 100 MB (right). Note that the first two rows correspond to searching on the uncompressed text. We report the running time required by **zsearch** and the relative time used by the other tools, i.e. the average running time used by **grep** when searching on 1 MB of JSON data is  $0.64 \cdot 5.33 = 3.4\text{ms}$ . For tools implementing the decompress and search approach we report the results obtained a) using the Linux command `taskset -c` to force decompressor and search to use a single CPU (without parenthesis) and b) when no restriction is imposed (in parenthesis).

**Analysis of the results.** It follows from Table 2 that **zsearch** outperforms **LZgrep** and **GNgrep**, the only two tools that perform the search on the compressed text without decompression. Besides, the compression algorithm on which **LZgrep** and **GNgrep** perform the search, LZW, exhibits much worse compression ratio than **repair**, as evidenced in Table 1.

<sup>15</sup> <http://storage.googleapis.com/books/ngrams/books/datasetsv2.html>

<sup>16</sup> <http://ita.ee.lbl.gov/html/contrib/NASA-HTTP.html>

<sup>17</sup> The files used for our experiments are available at TODO

<sup>18</sup> For **LZgrep** Table 2 shows the average running time required to search for the patterns in the list.

<sup>19</sup> <https://pevalme.github.io/zsearch/graphs/index.html>

Comparing **zearch**'s running time on different files we observe that its performance is related to the compression ratio achieved for the given file. This is to be expected since bigger compression ratio means there are less rules in the resulting grammar, each of which is processed by **zearch** exactly once. Hence, **zearch** appears as the fastest tool for searching on a compressed log file, being more than twice as fast as any other tool searching in compressed text and even outperforming **grep** and **rg**, which operate on the uncompressed text. Similarly, **zearch** yields its worst results when working on *CSV*.

Finally remark that, except for *CSV* files, **zearch** offers better performance than any other tool searching on compressed text using a single CPU and even when the decompress and search approach uses two processing units, it is outperformed by **zearch** in most of the experiments. We believe these results evidence the virtues of our approach, which allows us to keep our files effectively compressed (see Table 1) and yet be able to efficiently perform regular expression matching (see Figure 1) on them.

**Theoretical Complexity and Behavior in Practice** In the experiments **zearch** does not run into the worst case scenario, which corresponds to Algorithm 2 adding a transition for each variable and each pair of states (maintaining a complete graph at each step of the algorithm). Such situation implies that each string generated by a grammar variable, i.e. each repeated factor in the uncompressed text, labels a path in the automaton between each pair of states. However, the upper bound for the worst case time complexity of **zearch** can be reduced to  $\mathcal{O}((p - |\sigma|)\pi\pi_q + |\sigma|\pi)$  where  $\pi = \max_{w \in \Sigma^*} |\{\langle q_1, q_2 \rangle \in Q^2 \mid q_1 \xrightarrow{w} q_2\}|$  and  $\pi_q = \max_{w \in \Sigma^*, q_1 \in Q} |\{q_2 \in Q \mid q_1 \xrightarrow{w} q_2\}|$ . In our experiments we observe that  $\pi$  and  $\pi_q$  are much smaller than  $s^2$  and  $s$  respectively.

Let  $s_\tau$  be the number of transitions in the automaton labeled with  $\tau$ . With the data structures described in Section 4, the runtime complexity of processing the rule  $X \rightarrow \alpha\beta$  with Algorithm 3 is  $\mathcal{O}(s_\alpha \times s_\beta)$ . Similarly, processing each symbol  $(\sigma)_i$  of the axiom rule  $X \rightarrow \sigma$  requires  $\mathcal{O}(s_{(\sigma)_i})$  time. Table 3 shows, for the set of experiments used to generate Table 2, the average value of  $s_\alpha \times s_\beta$  and  $s_{(\sigma)_i}$  for each rule and axiom symbol of the grammar. Note that the given numbers are much smaller than the theoretical  $s^3$ .

File	$p -  \sigma $	$s_\alpha \times s_\beta$	$ \sigma $	$s_{(\sigma)_i}$	File	$p -  \sigma $	$s_\alpha \times s_\beta$	$ \sigma $	$s_{(\sigma)_i}$
Subtitles	34,056	8.86	32,426	1.55	Subtitles	2,646,418	3.37	2,367,739	1.32
CSV	142,012	12.18	24,992	0.83	CSV	9,518,105	1.91	803,548	0.51
JSON	35,282	8.36	27,013	1.06	JSON	1,929,180	3.19	1,428,333	0.77
Logs	28,476	6.41	12,329	0.54	Logs	2,345,505	2.58	455,132	0.46

**Table 3** Average number of operations performed by Algorithm 2 per grammar rule. The table on the left (resp. right) corresponds to input files of uncompressed size 1 MB (resp. 100MB).

## 6 Conclusions and Future Work

We have shown that the performance of regular expression matching on compressed text can be improved by working directly on the compressed file rather than decompressing it and searching on the output as it is generated. Furthermore we provide a tool, **zearch**, that despite being purely sequential outperforms the state of the art even when decompression and search are done simultaneously.

It is yet to be considered how the performance of **zearch** is affected by the choice of the grammar-based compression algorithm. By using different heuristics to build the grammar, the resulting SLP will have different properties, such as depth, width or length of the axiom

rule, which affect `zsearch`'s performance. Furthermore, there are grammar-based compression algorithms such as `Sequitur` that produce SLPs in which rules might have more than two symbols on the right hand side. It is worth considering whether adapting `zsearch` to work on such SLPs will have a positive impact on its performance.

Finally, remark that the current implementation is purely sequential. However, the algorithms involved allow for a conceptually simple parallelization since any set of rules such that the sets of symbols on the left and on the right hand side are disjoint can be processed simultaneously. Indeed, a theoretical result by Ullman et al. [21] on the parallelization of Datalog queries can be interpreted in our setting to show that the regular expression search on grammar-compressed text is in  $\mathcal{NC}^2$  when the automaton built from the expression is acyclic. Therefore it is worth considering the development of a parallel version of `zsearch`.

## References

- 1 W. Plandowski and W. Rytter. Complexity of language recognition problems for compressed words. In *Jewels are forever*, 1999.
- 2 N. Markey and Ph. Schnoebelen. A ptime-complete matching problem for slp-compressed words. *Information Processing Letters*, 2004.
- 3 A. Abboud, A. Backurs, K. Bringmann, and M. Künnemann. Fine-grained complexity of analyzing compressed data: Quantifying improvements over decompress-and-solve. *FOCS*, 2018.
- 4 G. Navarro and J. Tarhio. Lzgrep: a boyer-moore string matching tool for ziv-lempel compressed text. *Software: Practice and Experience*, 2005.
- 5 E. S. de Moura, G. Navarro, N. Ziviani, and R. A. Baeza-Yates. Direct pattern matching on compressed text. In *SPIRE*, 1998.
- 6 G. Navarro and V. Mäkinen. Compressed full-text indexes. *ACM Computing Surveys*, 2007.
- 7 P. Bille, R. Fagerberg, and I. L. Gørtz. Improved approximate string matching and regular expression matching on ziv-lempel compressed texts. *ACM Transactions on Algorithms*, 2009.
- 8 J. Kärkkäinen, G. Navarro, and E. Ukkonen. Approximate string matching on ziv-lempel compressed text. *Journal of Discrete Algorithms*, 2003.
- 9 T. Kida, M. Takeda, A. Shinohara, M. Miyazaki, and S. Arikawa. Multiple pattern matching in lzw compressed text. In *DCC*, 1998.
- 10 P. Gawrychowski. Simple and efficient lzw-compressed multiple pattern matching. *Journal of Discrete Algorithms*, 2014.
- 11 G. Navarro. Regular expression searching on compressed text. *Journal of Discrete Algorithms*, 2003.
- 12 M. Charikar, E. Lehman, Ding Liu, R. Panigrahy, M. Prabhakaran, A. Sahai, and A. Shelat. The smallest grammar problem. *IEEE Transactions on Information Theory*, 2005.
- 13 J. Ziv and A. Lempel. Compression of individual sequences via variable-rate coding. *IEEE Transactions on Information Theory*, 1978.
- 14 Terry A. Welch. A technique for high-performance data compression. *Computer*, 1984.
- 15 N. J. Larsson and A. Moffat. Off-line dictionary-based compression. *DCC*, 1999.
- 16 C. G. Nevill-Manning and I. H. Witten. Compression and explanation using hierarchical grammars. *The Computer Journal*, 1997.
- 17 D. HucKe, M. Lohrey, and C. Philipp Reh. The smallest grammar problem revisited. In *SPIRE*, 2016.
- 18 J. Esparza, P. Rossmanith, and S. Schwoon. A uniform framework for problems on context-free grammars. *Bulletin of the EATCS*, 2000.

- 430 **19** K. Thompson. Programming techniques: Regular expression search algorithm. *Commu-*  
431 *nications of the ACM*, 1968.
- 432 **20** P. Lison and J. Tiedemann. Opensubtitles2016: Extracting large parallel corpora from  
433 movie and TV subtitles. In *LREC*, 2016.
- 434 **21** J. D. Ullman and A. Van Gelder. Parallel complexity of logical query programs. *Algorith-*  
435 *mica*, 1988.

## A Deferred Proofs

### A.1 Proof of Theorem 1

The following definitions allow us to prove Theorem 1.

► **Definition 13** (Extended label). Let  $A = (Q, \Sigma, I, F, \delta)$  be an FSA,  $w \in \Sigma^+$  and  $q_1, q_2 \in Q$ . We say  $w$  is an *extended label* of a path from  $q_1$  to  $q_2$ , and write  $q_1 \rightsquigarrow^w q_2$ , iff there exists a run  $q_1 \xrightarrow{(w)_{i,j}} q_2$  with  $1 \leq i \leq j \leq |w|$  and  $(i = 1 \vee q_1 \in I) \wedge (j = |w| \vee q_2 \in F)$ . We refer to factor  $(w)_{i,j}$  as the *witness* of the extended label.

► **Definition 14** (Minimal witness). Let  $A = (Q, \Sigma, I, F, \delta)$  be an FSA,  $w \in \Sigma^+$  and  $q_1, q_2 \in Q$ . Let  $q_1 \rightsquigarrow^w q_2$  with witness  $w' \in \Sigma^+$ . We say witness  $w'$  is a *minimal witness* if no factor of  $w'$  is a witness of  $q_1 \rightsquigarrow^w q_2$ .

Intuitively,  $q_1 \rightsquigarrow^w q_2$  indicates that if transitions  $\{(q, a, q) \mid q \in (I \cup F), a \in \Sigma\}$  were added to  $\delta$  then there would be a run  $q_1 \xrightarrow{w} q_2$ . Proofs of Lemmas 15 and 16 rely on this concept to describe the set of transitions that Algorithm 2 adds to the automaton.

► **Lemma 15.** Let  $P = (V, \Sigma, \mathcal{R})$  be an SLP and  $A = (Q, \Sigma, I, F, \delta)$  an FSA with  $I \cap F = \emptyset$ . Let  $X \in V$ ,  $w \in \Sigma^*$  with  $X \Rightarrow^* w$  and  $q_1, q_2 \in Q$ . After  $\ell$  iterations of the outermost “for” loop of Algorithm 1, if  $X \in \{X_1, \dots, X_\ell\} \subseteq V$  then

$$(q_1, X, q_2) \in \delta \implies q_1 \rightsquigarrow^w q_2.$$

**Proof.**

**Base case:**  $\ell = 1$ .

By definition of SLP,  $X_1 \rightarrow ab$  with  $ab \in \Sigma$ . Let  $w = ab$  so  $X_1 \Rightarrow^* w$ . Since  $I \cap F = \emptyset$ , there are three possible scenarios for which transition  $(q_1, X_1, q_2)$  is added:

- $(q_1, a, q'), (q', b, q_2) \in \delta$ . Then  $q_1 \xrightarrow{ab} q_2$  so  $ab$  is a witness of  $q_1 \rightsquigarrow^w q_2$ .
- $q' = q_2 \in F$  and  $(q_1, a, q') \in \delta$ . Then  $q_1 \xrightarrow{a} q'$  so  $a$  is a witness of  $q_1 \rightsquigarrow^w q_2$ .
- $q_1 = q' \in I$  and  $(q', b, q_2) \in \delta$ . Then  $q' \xrightarrow{b} q_2$  so  $b$  is a witness of  $q_1 \rightsquigarrow^w q_2$ .

**Inductive step:** Assume the lemma holds for some  $\ell < |\mathcal{R}|$  and consider rule  $X_{\ell+1} \rightarrow \alpha\beta$  with  $X_{\ell+1} \Rightarrow^* w$  for some  $w \in \Sigma^*$ . By definition of SLP  $\alpha, \beta \in (\Sigma \cup \{X_1, X_2, \dots, X_\ell\})$  and  $X \Rightarrow \alpha\beta \Rightarrow^* u\beta \Rightarrow^* uv = w$ . Again, there are three possible scenarios for which transition  $(q_1, X_{\ell+1}, q_2)$  is added:

- $(q_1, \alpha, q'), (q', \beta, q_2) \in \delta$ . By hypothesis,  $q_1 \rightsquigarrow^u q'$  and  $q' \rightsquigarrow^v q_2$  so, by Definition 13,  $q_1 \rightsquigarrow^w q_2$ .
- $q' = q_2 \in F$  and  $(q_1, \alpha, q') \in \delta$ . By hypothesis,  $q_1 \rightsquigarrow^u q_2$  so, by Definition 13,  $q_1 \rightsquigarrow^w q_2$ .
- $q_1 = q' \in I$  and  $(q', \beta, q_2) \in \delta$ . By hypothesis,  $q_1 \rightsquigarrow^v q_2$  so, by Definition 13,  $q_1 \rightsquigarrow^w q_2$ . ◀

► **Lemma 16.** Let  $P = (V, \Sigma, \mathcal{R})$  be an SLP and  $A = (Q, \Sigma, I, F, \delta)$  an FSA with  $I \cap F = \emptyset$ . Let  $X \in V$ ,  $w \in \Sigma^*$  with  $X \Rightarrow^* w$  and  $q_1, q_2 \in Q$ . After  $\ell$  iterations of the outermost “for” loop of Algorithm 1, if  $X \in \{X_1, \dots, X_\ell\} \subseteq V$  then

$$q_1 \rightsquigarrow^w q_2 \implies (q_1, X, q_2) \in \delta.$$

**Proof.**

**Base case:**  $\ell = 1$ .

By definition of SLP,  $X_1 \rightarrow ab$  with  $a, b \in \Sigma$  and, since  $I \cap F = \emptyset$ , there are three possibilities for  $w' \in \Sigma^+$ , the *minimal witness* of  $q_1 \rightsquigarrow^{ab} q_2$ :

- $w' = ab$ . Then  $(q_1, a, q_3), (q_3, b, q_2) \in \delta$  so  $(q_1, X_1, q_2)$  is added to  $\delta$  with  $q' = q_3$ .

- 476 ■  $w' = a$ . Then  $(q_1, a, q_2) \in \delta$  and  $q_2 \in F$  so  $(q_1, X_1, q_2)$  is added to  $\delta$  with  $q' = q_2 \in F$ .
- 477 ■  $w' = b$ . Then  $(q_1, b, q_2) \in \delta$  and  $q_1 \in I$  so  $(q_1, X_1, q_2)$  is added to  $\delta$  with  $q_1 = q' \in I$ .

478 **Inductive step:** Assume the lemma holds for some  $\ell < |\mathcal{R}|$  and consider rule  $X_{\ell+1} \rightarrow \alpha\beta$   
 479 with  $X_{\ell+1} \Rightarrow^* w$  for some  $w \in \Sigma^*$ . By definition of SLP  $\alpha, \beta \in (\Sigma \cup \{X_1, X_2, \dots, X_\ell\})$  and  
 480  $\alpha \Rightarrow^* u, \beta \Rightarrow^* v$  with  $w = u \cdot v$ . Again, there are three possibilities for  $w' \in \Sigma^+$ , the *minimal*  
 481 *witness* of  $q_1 \rightsquigarrow^w q_2$ .

- 482 ■  $w'$  is the concatenation of a nonempty suffix of  $u$  and a nonempty prefix of  $v$ . By  
 483 Definition 13  $q_1 \rightsquigarrow^u q_3$  and  $q_3 \rightsquigarrow^v q_2$  for some  $q_3 \in Q$ . By hypothesis  $(q_1, \alpha, q_3), (q_3, \beta, q_2) \in \delta$   
 484 and therefore  $(q_1, X_{\ell+1}, q_2)$  is added to  $\delta$  when  $q' = q_3$ .
- 485 ■  $w'$  is a factor of  $u$ . By Definition 13  $q_2 \in F$  and  $q_1 \rightsquigarrow^u q_2$ . By hypothesis  $(q_1, \alpha, q_2) \in \delta$  so  
 486 transition  $(q_1, X_{\ell+1}, q_2)$  is added to  $\delta$  when  $q' = q_2$ .
- 487 ■  $w'$  is a factor of  $v$ . Definition 13  $q_1 \in I$  and  $q_1 \rightsquigarrow^v q_2$ . By hypothesis  $(q_1, \beta, q_2) \in \delta$  so  
 488 transition  $(q_1, X_{\ell+1}, q_2)$  is added to  $\delta$  with  $q_1 = q'$ . ◀

489 **Proof of Theorem 1.** Assume  $(q_0, X, q_f) \in \delta$  with  $q_0 \in I$  and  $q_f \in F$ . Then by Lemma 15,  
 490  $q_0 \rightsquigarrow^w q_f$  and from Definition 13 it follows  $w \in (\Sigma^* \cdot \mathcal{L}(A) \cdot \Sigma^*)$ .

491 Now assume  $w \in (\Sigma^* \cdot \mathcal{L}(A) \cdot \Sigma^*)$ . There is a factor  $w'$  for which  $q_0 \xrightarrow{w'} q_f$  with  $q_0 \in I$  and  $q_f \in F$   
 492 so, by Definition 13,  $q_0 \rightsquigarrow^w q_f$ . By Lemma 16,  $(q_0, X, q_f)$  is added to  $\delta$  by Algorithm 1. ◀

## 493 A.2 Proof of Lemma 5

- 494 1. It follows from Definition 3 since  $(w)_{i,j} = (u)_{i,j}$  for all  $1 \leq i \leq j < |u|$ .
- 495 2. It follows from Definition 3 since  $(w)_{i+|u|,j+|u|} = (v)_{i,j}$  for all  $1 < i \leq j \leq |v|$ . ◀

## 496 A.3 Proof of Lemma 6

497 Note that  $(w)_i = (u)_i$  for  $1 \leq i \leq |u|$  and  $(w)_{i+|u|} = (v)_i$  for  $1 \leq i \leq |v|$ . Therefore,  
 498 it follows from Definition 3 that if there exists  $(i, j) \in \mathcal{ML}_X$  with  $i \leq |u| < j$  then  
 499  $(i, j) = (\ell_X, r_X)$ . Next we prove that if any of the properties holds then  $(\ell_X, r_X) \in \mathcal{ML}_X$ .  
 500 Let  $\widehat{\mathcal{L}}(A) = \widehat{\Sigma}^* \cdot \mathcal{L}(A) \cdot \widehat{\Sigma}^*$ .

- 501 1.  $(\ell_X, |u|) \in \mathcal{ML}_\alpha$  means that  $(w)_{\ell_X, |u|} \in \widehat{\mathcal{L}}(A)$ . Since  $(w)_{\ell_X, |u|}$  is a factor of  $(w)_{\ell_X, r_X} \in \widehat{\Sigma}^+$   
 502 it follows that  $(w)_{\ell_X, r_X} \in \widehat{\mathcal{L}}(A)$  and, therefore,  $(\ell_X, r_X) \in \mathcal{ML}_X$ .
- 503 2.  $(1, r_X - |u|) \in \mathcal{ML}_\beta$  means that  $(w)_{1+|u|, r_X} \in \widehat{\mathcal{L}}(A)$ . Since  $(w)_{1+|u|, r_X}$  is a factor of  
 504  $(w)_{\ell_X, r_X} \in \widehat{\Sigma}^+$  it follows that  $(w)_{\ell_X, r_X} \in \widehat{\mathcal{L}}(A)$  and, therefore,  $(\ell_X, r_X) \in \mathcal{ML}_X$ .
- 505 3. By Lemma 15  $q_0 \rightsquigarrow^u q'$  and  $q' \rightsquigarrow^v q_f$  with *witnesses*  $u'$  and  $v'$  respectively. Since  $\prec \notin \widehat{\Sigma}$  and  
 506  $q' \notin (I \cup F)$  then  $u' \cdot v'$  is a factor of  $(w)_{\ell_X, r_X} \in \widehat{\mathcal{L}}(A)$  and, therefore,  $(\ell_X, r_X) \in \mathcal{ML}_X$ .

507 We conclude by proving the other side of the implication: if  $(\ell_X, r_X) \in \mathcal{ML}_X$  then at least  
 508 one of the properties a)-c) holds. By Definition 3,  $(\ell_X, r_X) \in \mathcal{ML}_X$  means  $(w)_{\ell_X, r_X} \in \widehat{\mathcal{L}}(A)$   
 509 and, therefore,  $q_0 \rightsquigarrow^w q_f$  with  $q_0 \in I$  and  $q_f \in F$ . There are three possibilities for  $w'$ , the  
 510 *minimal witness* of  $q_0 \rightsquigarrow^w q_f$ :

- 511 ■  $w'$  is a factor of  $u$ . Then  $q_0 \rightsquigarrow^u q_f$  and, therefore  $(w)_{\ell_X, |u|} \in \widehat{\mathcal{L}}(A)$  so property a) holds.
- 512 ■  $w'$  is a factor of  $v$ . Then  $q_0 \rightsquigarrow^v q_f$  and, therefore  $(w)_{|u|, r_X} \in \widehat{\mathcal{L}}(A)$  so property b) holds.
- 513 ■ Otherwise  $w'$  is the concatenation of a nonempty suffix of  $u$  and a nonempty prefix of  $v$ .  
 514 Then  $q_0 \rightsquigarrow^u q', q' \rightsquigarrow^v q_f$  for some  $q' \in (Q \setminus (I \cup F))$  and, by Lemma 16, property c) holds. ◀



#### 515 A.4 Proof of Lemma 10

516 Let  $X \Rightarrow^* w$ ,  $\alpha \Rightarrow^* u$ ,  $\beta \Rightarrow^* v$ . Next we show that, for each element in  $\mathcal{C}_X$ , the value given  
517 by Definition 8 coincides with the one computed using the formulas given by Lemma 10.

- 518 1.  $\exists k (w)_k = \prec \triangleright$  iff  $(\exists k (u)_k = \prec \triangleright) \vee (\exists k (v)_k = \prec \triangleright)$ . Therefore  $\text{NL}(X) = \text{NL}(\alpha) \vee \text{NL}(\beta)$ .  
519 2. We provide the proof for  $\text{L}(X)$  but not for  $\text{R}(X)$  since they are conceptually identical.  
520  $\text{L}(X) := \exists i$  s.t.  $(0, i) \in \mathcal{ML}_\tau$ . If  $(1, j) \in \mathcal{ML}_X$  with  $j < |u|$  then  $\text{NL}(\alpha) = \text{true}$  and,  
521 by Lemma 5,  $(1, j) \in \mathcal{ML}_X$  iff  $(1, j) \in \mathcal{ML}_\alpha$ , i.e.  $\text{L}(\alpha) = \text{true}$ . Otherwise  $j = r_X$ ,  
522  $\text{NL}(\alpha) = \text{false}$  and, by Lemma 6,  $(1, r_X) \in \mathcal{ML}_X$  iff one of the following holds:
- 523 a.  $(1, |u|) \in \mathcal{ML}_\alpha$  i.e.  $\text{L}(\alpha) = \text{true}$ .
  - 524 b.  $(1, r_X - |u|) \in \mathcal{ML}_\beta$  i.e.  $\text{L}(\beta) = \text{true}$ .
  - 525 c. Algorithm 1 adds  $(q_0, \alpha, q')$ ,  $(q', \beta, q_f)$  to  $\delta$  with  $q_0 \in I$ ,  $q' \notin (I \cup F)$  and  $q_f \in F$  i.e.  
526 function COUNT is invoked with  $m = \text{true}$ .

527 Therefore  $\text{L}(X) = \begin{cases} \text{L}(\alpha) \vee \text{L}(\beta) \vee m & \text{if } \neg \text{NL}(\alpha) \\ \text{L}(\alpha) & \text{otherwise} \end{cases}$

- 528 3. Write  $\mathcal{ML}_X$  as the union of three disjoint sets as follows:

529  $\mathcal{ML}_X = \{(i, j) \in \mathcal{ML}_X \mid i, j < |u|\} \cup \{(i, j) \in \mathcal{ML}_X \mid |u| < i, j\} \cup \{(i, j) \in \mathcal{ML}_X \mid i \leq |u| \leq j\}$ .  
530 From Lemma 5 and Definition 8 it follows  $\text{M}(X) = \text{M}(\alpha) + \text{M}(\beta) + |\{(i, j) \in \mathcal{ML}_{\mathcal{C}_X} \mid i \leq |u| < j\}|$   
531 with  $\{(i, j) \in \mathcal{ML}_{\mathcal{C}_X} \mid i \leq |u| \leq j\} \neq \emptyset$  iff  $(\ell_X, r_X) \in \mathcal{ML}_X$  with  $\ell_X \neq 1$  and  $r_X \neq |w|$   
532 which, by Lemma 6, occurs iff one of the following holds and  $\text{NL}(\alpha) = \text{NL}(\beta) = \text{true}$ :

- 533 a.  $(\ell_X, |u|) \in \mathcal{ML}_\alpha$  i.e.  $\text{L}(\alpha) = \text{true}$ .
- 534 b.  $(\ell_X, r_X - |v|) \in \mathcal{ML}_\beta$  i.e.  $\text{R}(\beta) = \text{true}$ .
- 535 c. Algorithm 1 adds  $(q_0, \alpha, q')$ ,  $(q', \beta, q_f)$  to  $\delta$  with  $q_0 \in I$ ,  $q' \notin (I \cup F)$  and  $q_f \in F$  i.e.  
536 function COUNT is invoked with  $m = \text{true}$ .

537 Therefore  $\text{M}(X) = \text{M}(\alpha) + \text{M}(\beta) + \begin{cases} 1 & \text{if } (\text{NL}(\alpha) \vee \text{NL}(\beta)) \wedge (\text{R}(\alpha) \vee \text{L}(\beta) \vee m) \\ 0 & \text{otherwise} \end{cases}$  ◀

## 538 B Algorithms

---

**Algorithm 3** Extension of Algorithm 2 to handle ESLP-compressed text

---

**Input:** An ESLP  $P = (V, \Sigma, \mathcal{R})$  and an FSA  $A = (Q, \Sigma, I, F, \delta)$

**Output:** Number of lines in  $\mathcal{L}(P)$  containing a factor matching  $A$

---

```

1: function MAIN
2:   INIT_AUTOMATON( )
3:   for each  $\ell = 1, 2, \dots, |V| - 1$  do
4:     execute lines [13–19] of Algorithm 2
5:   let  $(X_{|V|} \rightarrow \sigma) \in \mathcal{R}$ 
6:    $\mathcal{C}_{X_{|V|}} := (\text{NL}(X_{|V|}) := 0, \text{L}(X_{|V|}) := 0, \text{R}(X_{|V|}) := 0, \text{M}(X_{|V|}) := 0)$ 
7:   active_states :=  $\emptyset$ 
8:   for each  $\lambda = 1, \dots, |\sigma|$  do
9:      $\mathcal{C}_{\text{tmp}} := \mathcal{C}_{X_{|V|}}$ 
10:    new_match := false
11:    active_states :=  $I \cup \{q_2 \mid (q_1, (\sigma)_\lambda, q_2) \in \delta, q_1 \in \text{active\_states}\}$ 
12:    for each  $(q_1, (\sigma)_\lambda, q_2) \in \delta$  s.t.  $q_1 \in \text{active\_states}$  do
13:      new_match := new_match  $\vee (q_2 \in F \wedge q_1 \notin (I \cup F))$ 
14:    COUNT( $\mathcal{C}_{X_{|V|}}, \mathcal{C}_{\text{tmp}}, \mathcal{C}_{(\sigma)_\lambda}, \text{new\_match}$ )
15:  return  $(\text{NL}(X_{|V|}) ? \text{M}(X_{|V|}) + \text{L}(X_{|V|}) + \text{R}(X_{|V|}) : \text{L}(X_{|V|}))$ 

```

---