

# 1 Counting on General Run-Length Grammars

2 **Gonzalo Navarro** ✉

3 Center for Biotechnology and Bioengineering (CeBiB)

4 Department of Computer Science, University of Chile, Chile

5 **Alejandro Pacheco** ✉

6 Center for Biotechnology and Bioengineering (CeBiB)

7 Department of Computer Science, University of Chile, Chile

---

## 8 — Abstract —

9 We introduce a data structure for counting pattern occurrences in texts compressed with any  
10 run-length context-free grammar. Our structure uses space proportional to the grammar size and  
11 counts the occurrences of a pattern of length  $m$  in a text of length  $n$  in time  $O(m \log^{2+\epsilon} n)$ , for  
12 any constant  $\epsilon > 0$  chosen at indexing time. This is the first solution to an open problem posed by  
13 Christiansen et al. [ACM TALG 2020] and enhances our abilities for computation over compressed  
14 data; we give an example application.

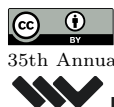
15 **2012 ACM Subject Classification** Theory of computation → Data structures design and analysis

16 **Keywords and phrases** Grammar-based indexing; Run-length context-free grammars, Counting  
17 pattern occurrences; Periods in strings.

18 **Digital Object Identifier** 10.4230/LIPIcs.CPM.2025.

19 **Funding** *Gonzalo Navarro*: Funded by Basal Funds FB0001, Mideplan, Chile, and Fondecyt Grant  
20 1-230755, Chile.

21 *Alejandro Pacheco*: Funded by Basal Funds FB0001, Mideplan, Chile, Fondecyt Grant 1-230755,  
22 Chile, and ANID/Scholarship Program/DOCTORADO BECAS CHILE/2018-21180760.



23 **1** Introduction

24 Context-free grammars (CFGs) have proven to be an elegant and efficient model for data  
 25 compression. The idea of grammar-based compression [47, 26] is, given a text  $T[1..n]$ ,  
 26 to construct a context-free grammar  $G$  of size  $g$  that only generates  $T$ . One can then  
 27 store  $G$  instead of  $T$ , which achieves compression if  $g \ll n$ . Compared to more powerful  
 28 compression methods like Lempel-Ziv [32], grammar compression offers efficient direct access  
 29 to arbitrary snippets of  $T$  without the need of full decompression [45, 2]. This has been  
 30 extended to offering indexed searches (i.e., in time  $o(n)$ ) for the occurrences of string patterns  
 31 in  $T$  [7, 14, 9, 6, 36], as well as more complex computations over the compressed sequence  
 32 [29, 19, 16, 17, 37, 25]. Since finding the smallest grammar  $G$  representing a given text  $T$  is  
 33 NP-hard [45, 4], many algorithms have been proposed to find small grammars for a given  
 34 text [31, 45, 42, 46, 33, 20, 21]. Grammar compression is particularly effective when handling  
 35 repetitive texts; indeed, the size  $g^*$  of the smallest grammar representing  $T$  is used as a  
 36 measure of its repetitiveness [35].

37 Nishimoto et al. [43] proposed enhancing CFGs with “run-length rules” to improve the  
 38 compression of repetitive strings. These run-length rules have the form  $A \rightarrow B^s$ , where  $B$  is  
 39 a terminal or a non-terminal symbol and  $s \geq 2$  is an integer. CFGs that may use run-length  
 40 rules are called run-length context-free grammars (RLCFGs). Because CFGs are RLCFGs,  
 41 the size  $g_{rl}^*$  of the smallest RLCFG generating  $T$  always satisfies  $g_{rl}^* \leq g^*$ , and it can be  
 42  $g_{rl}^* = o(g^*)$  in text families as simple as  $T = a^n$ , where  $g_{rl}^* = O(1)$  and  $g^* = \Theta(\log n)$ .

43 The use of run-length rules has become essential to produce grammars with size guarantees  
 44 and convenient regularities that speed up indexed searches and other computations [29, 19,  
 45 16, 6, 25, 27]. The progress made in indexing texts with CFGs has been extended to RLCFGs,  
 46 reaching the same status in most cases. These functionalities include extracting substrings,  
 47 computing substring summaries, and locating all the occurrences of a pattern string [6,  
 48 App. A]. It has also been shown that RLCFGs can be balanced [38] in the same way as CFGs  
 49 [17], which simplifies many compressed computations on RLCFGs.

50 Interestingly, *counting*, that is, determining how many times a pattern occurs in the text  
 51 without spending the time to list those occurrences, can be done efficiently on CFGs, but  
 52 not so far on RLCFGs. Counting is useful in various fields, such as pattern discovery and  
 53 ranked retrieval, for example to help determine the frequency or relevance of a pattern in  
 54 the texts of a collection [34].

55 Navarro [40] showed how to count the occurrences of a pattern  $P[1..m]$  in  $T[1..n]$  in  
 56  $O(m^2 + m \log^{2+\epsilon} n)$  time using  $O(g)$  space if a CFG of size  $g$  represents  $T$ , for any constant  
 57  $\epsilon > 0$  chosen at indexing time. Christiansen et al. improved this time to  $O(m \log^{2+\epsilon} n)$  by  
 58 using more recent underlying data structures for tries. Christiansen et al. [6] and Kociumaka  
 59 et al. [27] extended the result to *particular* RLCFGs, even achieving optimal  $O(m)$  time by  
 60 using additional space, but could not extend their mechanism to general RLCFGs.

61 In this paper we give the first solution to this open problem, by introducing an index  
 62 that counts the occurrences of a pattern  $P[1..m]$  in a text  $T[1..n]$  represented by a RLCFG  
 63 of size  $g_{rl}$ . Our index uses  $O(g_{rl})$  space and answers queries in time  $O(m \log^{2+\epsilon} n)$  for any  
 64 constant  $\epsilon > 0$  chosen at indexing time. This is the same time complexity that holds for  
 65 CFGs, which puts on par our capabilities to handle RLCFGs and CFGs on all the considered  
 66 functionalities. As an example of our new capabilities, we show how a recent result on finding  
 67 the maximal exact matches of  $P$  using CFGs [41] can now run on RLCFGs.

68 While our solution builds on the ideas developed for CFGs and particular RLCFGs  
 69 [40, 6, 27], arbitrary RLCFGs lack crucial structure that holds in those particular cases,

70 namely that if there exists a run-length rule  $A \rightarrow B^s$ , then the period [10] of the string  
 71 represented by  $A$  is the length of that of  $B$ . We show, however, that the general case still  
 72 retains some structure relating the shortest periods of  $P$  and the string represented by  $A$ .  
 73 We exploit this relation to develop a solution that, while considerably more complex than  
 74 that for those particular cases, retains the same theoretical guarantees obtained for CFGs.

## 75 2 Basic Concepts

### 76 2.1 Strings

77 A *string*  $S[1..n] = S[1] \cdot S[2] \cdots S[n]$  is a sequence of symbols, where each symbol belongs  
 78 to a finite ordered set of integers called an *alphabet*  $\Sigma = \{1, 2, \dots, \sigma\}$ . The *length* of  $S$  is  
 79 denoted by  $|S| = n$ . We denote with  $\varepsilon$  the empty string, where  $|\varepsilon| = 0$ . A *substring* of  $S$  is  
 80  $S[i..j] = S[i] \cdot S[i+1] \cdots S[j]$  (which is  $\varepsilon$  if  $i > j$ ). A *prefix* (*suffix*) is a substring of the  
 81 form  $S[1..j] = S[1] \cdots S[j]$  ( $S[j..n] = S[j] \cdots S[n]$ ); we also say that  $S[1..j]$  ( $S[j..n]$ ) *prefixes* (*suffixes*)  
 82  $S$ . We write  $S \sqsubseteq S'$  if  $S$  prefixes  $S'$ , and  $S \sqsubset S'$  if in addition  $S \neq S'$  ( $S$  strictly prefixes  $S'$ ).

83 We denote with  $S \cdot S'$  the *concatenation* of  $S$  and  $S'$ . A *power*  $t \in \mathbb{N}$  of a string  $S$ , written  
 84  $S^t$ , is the concatenation of  $t$  copies of  $S$ . The *reverse* string of  $S[1..n] = S[1] \cdot S[2] \cdots S[n]$   
 85 refers to  $S[1..n]^{\text{rev}} = S[n] \cdot S[n-1] \cdots S[1]$ . We also use the term *text* to refer to a string.

### 86 2.2 Periods of strings

87 Periods of strings [10] are crucial in this paper. We recall their definition(s) and a key  
 88 property, the renowned Periodicity Lemma.

- 89 ► **Definition 1.** A string  $S[1..n]$  has a period  $1 \leq p \leq n$  if, equivalently,  
 90 1. it consists of  $\lfloor n/p \rfloor$  consecutive copies of  $S[1..p]$  plus a (possibly empty) prefix of  $S[1..p]$ ,  
 91 that is,  $S = (S[1..p]^{\lfloor n/p \rfloor})[1..n]$ ; or  
 92 2.  $S[1..n-p] = S[p+1..n]$ ; or  
 93 3.  $S[i+p] = S[i]$  for all  $1 \leq i \leq n-p$ .

94 We also say that  $p$  is a period of  $S$ . We define  $p(S)$  as the shortest period of  $S$  and say  $S$  is  
 95 periodic if  $p(S) \leq n/2$ .

96 ► **Lemma 2** ([13]). If  $p$  and  $p'$  are periods of  $S$  and  $|S| \geq p + p' - \gcd(p, p')$ , then  $\gcd(p, p')$   
 97 is a period of  $S$ . Thus,  $p(S)$  divides all other periods  $p \leq |S|/2$  of  $S$ .

### 98 2.3 Karp-Rabin signatures

99 Karp-Rabin [23] fingerprinting assigns a function  $k(S) = (\sum_{i=1}^m S[i] \cdot c^{i-1}) \bmod \mu$  to the  
 100 string  $S[1..m]$ , where  $c$  is a suitable integer and  $\mu$  a prime number. Bille et al. [3] showed  
 101 how to build, in  $O(n \log n)$  expected time, a *Karp-Rabin signature*  $\kappa(S)$  built from a pair of  
 102 Karp-Rabin functions, which has no collisions between substrings  $S$  of  $T[1..n]$ . We always  
 103 assume those kind of signatures in this paper.

104 A well-known property is that we can compute the signatures of all the prefixes  $S[1..j] \sqsubseteq S$   
 105 in time  $O(m)$ , and then obtain any  $\kappa(S[i..j])$  in constant time by using arithmetic operations.

### 106 2.4 Range summary queries on grids

107 A discrete grid of  $r$  rows and  $c$  columns stores points at integer coordinates  $(x, y)$ , with  
 108  $1 \leq x \leq c$  and  $1 \leq y \leq r$ . Grids with  $m$  points can be stored in  $O(m)$  space, so that some  
 109 *summary* queries are performed on *orthogonal ranges* of the grid. In particular, one can

## XX:4 Counting on General Run-Length Grammars

110 associate an integer with each point, and then, given an orthogonal range  $[x_1, x_2] \times [y_1, y_2]$ ,  
111 compute the *sum* of all the integers associated with the points in that range. Chazelle [5]  
112 showed how to run that query in time  $O(\log^{2+\epsilon} m)$ , for any constant  $\epsilon > 0$ , in  $O(m)$  space.

### 113 2.5 Grammar compression and parse trees

114 A *context-free grammar* (CFG)  $G = (V, \Sigma, R, S)$  is a language generation model consisting of  
115 a finite set of nonterminal symbols  $V$  and a finite set of terminal symbols  $\Sigma$ , disjoint from  $V$ .  
116 The set  $R$  contains a finite set of production rules  $A \rightarrow \alpha$ , where  $A$  is a nonterminal symbol  
117 and  $\alpha$  is a string of terminal and nonterminal symbols. The language generation process  
118 starts from a sequence formed by just the nonterminal  $S \in V$  and, iteratively, chooses a rule  
119  $A \rightarrow \alpha$  and replaces an occurrence of  $A$  in the sequence by  $\alpha$ , until the sequence contains  
120 only terminals. The *size* of the grammar,  $g = |G|$ , is the sum of the lengths of the right-hand  
121 sides of the rules,  $g = \sum_{A \rightarrow \alpha \in R} |\alpha|$ . Given a string  $T$ , we can build a CFG  $G$  that generates  
122 only  $T$ . Then, especially if  $T$  is repetitive,  $G$  is a compressed representation of  $T$ . The  
123 *expansion*  $\text{exp}(A)$  of a nonterminal  $A$  is the string generated by  $A$ , for instance  $\text{exp}(S) = T$ ;  
124 for terminals  $a$  we also say  $\text{exp}(a) = a$ . We use  $|A| = |\text{exp}(A)|$  and  $p(A) = p(\text{exp}(A))$ .

125 The *parse tree* of a grammar is an ordinal labeled tree where the root is labeled with  
126 the initial symbol  $S$ , the leaves are labeled with terminal symbols, and internal nodes are  
127 labeled with nonterminals. If  $A \rightarrow \alpha_1 \cdots \alpha_t$ , with  $\alpha_i \in V \cup \Sigma$ , then a node  $v$  labeled  $A$  has  $t$   
128 children labeled, left to right,  $\alpha_1, \dots, \alpha_t$ . A more compact version of the parse tree is the  
129 *grammar tree*, which is obtained by pruning the parse tree such that only one internal node  
130 labeled  $A$  is kept for each nonterminal  $A$ , while the rest become leaves. Unlike the parse  
131 tree, the grammar tree of  $G$  has only  $g + 1$  nodes. Consequently, the text  $T$  can be divided  
132 into at most  $g$  substrings, called *phrases*, each being the expansion of a grammar tree leaf.  
133 The starting phrase positions constitute a *string attractor* of the text [24]. Therefore, all text  
134 substrings of length more than 1 have at least one occurrence that crosses a phrase boundary.

### 135 2.6 Run-length grammars

136 *Run-length CFGs* (RLCFGs) [43] extend CFGs by allowing in  $R$  rules of the form  $A \rightarrow \beta^s$ ,  
137 where  $s \geq 2$  is an integer and  $\beta$  is a string of terminals and nonterminals. These rules are  
138 equivalent to rules  $A \rightarrow \beta \cdots \beta$  with  $s$  repetitions of  $\beta$ . However, the length of the right-hand  
139 side of the rule  $A$  is defined as  $|\beta| + 1$ , not  $s \cdot |\beta|$ . To simplify, we will only allow run-length  
140 rules of the form  $A \rightarrow B^s$ , where  $B$  is a single terminal or nonterminal; this does not increase  
141 the asymptotic grammar size because we can rewrite  $A \rightarrow B^s$  and  $B \rightarrow \beta$  for a fresh  $B$ .

142 RLCFGs are never larger than general CFGs, and they can be asymptotically smaller.  
143 For example, the size  $g_{rl}^*$  of the smallest RLCFG that generates  $T$  is in  $O(\delta \log \frac{n \log |\Sigma|}{\delta \log n})$ ,  
144 where  $\delta$  is a measure of repetitiveness based on substring complexity [44, 28], but such a  
145 bound does not always hold for the size  $g^*$  of the smallest grammar. The maximum stretch  
146 between  $g^*$  and  $g_{rl}^*$  is  $O(\log n)$ , as we can replace each rule  $A \rightarrow B^s$  by  $O(\log s)$  CFG rules.

147 We denote the size of an RLCFG  $G$  as  $g_{rl} = |G|$ . To maintain the invariant that the  
148 grammar tree has  $g_{rl} + 1$  nodes, we represent rules  $A \rightarrow B^s$  as a node labeled  $A$  with two  
149 children: the first is  $B$  and the second is a special leaf  $B^{[s-1]}$ , denoting  $s - 1$  repetitions of  $B$ .

## 150 3 Grammar Indexing for Locating

151 A *grammar index* represents a text  $T[1..n]$  using a grammar  $G$  that generates only  $T$ . As  
152 opposed to mere compression, the index supports three primary pattern-matching queries:

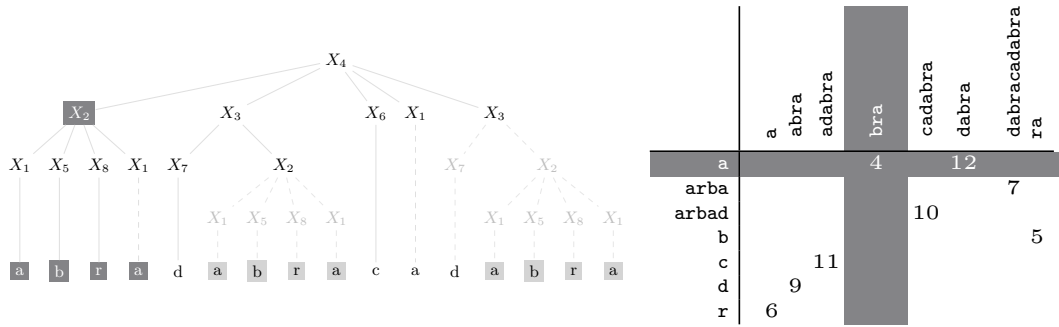


Figure 1 On the left, a grammar tree for  $T = abracadabra$  (with straight solid edges), so  $\text{exp}(X_4) = T$ . Dashed edges were removed from the parse tree. The only primary occurrence of  $P = abra$  in  $T$  is marked with dark gray on the bottom; the secondary ones are in light gray. On the right, the grid used for searching primary occurrences. Gray stripes indicate the search ranges corresponding to the partition  $P = R | Q$ , where  $R = a$  and  $Q = bra$ . The value 4 stored in the resulting cell is the preorder of the child  $X_5$  of the locus node  $X_2$  where  $Q$  starts.

153 *locate* (returning all positions of a pattern in the text), *count* (returning the number of times  
 154 a pattern appears in the text), and *extract* (extracting any desired substring of  $T$ ). In order  
 155 to locate, grammar indexes identify “initial” pattern occurrences and then track their “copies”  
 156 throughout the text. The former are the *primary occurrences*, defined as those that cross  
 157 phrase boundaries, and the latter are the *secondary occurrences*, which are confined to a  
 158 single phrase. This approach [22] forms the basis of most grammar indexes [7, 8, 9] and  
 159 related ones [14, 30, 11, 15, 12, 1, 39, 48], which first locate the primary occurrences and  
 160 then derive their secondary occurrences through the grammar tree.

161 As mentioned in Section 2.5, the grammar tree leaves cut the text into phrases. In order  
 162 to report each primary occurrence of a pattern  $P[1..m]$  exactly once, let  $v$  be the lowest  
 163 common ancestor of the first and last leaves the occurrence spans;  $v$  is called the *locus*  
 164 *node* of the occurrence. Let  $v$  have  $t$  children and the first leaf that covers the occurrence  
 165 descend from the  $i$ th child of  $v$ . If  $v$  represents  $A \rightarrow \alpha_1 \cdots \alpha_t$ , it follows that  $\text{exp}(\alpha_i)$  finishes  
 166 with a pattern prefix  $R = P[1..q]$  and that  $\text{exp}(\alpha_{i+1}) \cdots \text{exp}(\alpha_t)$  starts with the suffix  
 167  $Q = P[q + 1..m]$ . We will denote such *cuts* as  $P = R | Q$ . The alignment of  $R | Q$  within  
 168  $\text{exp}(\alpha_i) | \text{exp}(\alpha_{i+1}) \cdots \text{exp}(\alpha_t)$  is the only possible one for that primary occurrence.

169 Following the original scheme [22], grammar indexing builds two sets of strings,  $\mathcal{X}$  and  $\mathcal{Y}$ ,  
 170 to find primary occurrences [7, 8, 9]. For each grammar rule  $A \rightarrow \alpha_1 \cdots \alpha_t$ , the set  $\mathcal{X}$  contains  
 171 all the reverse expansions of the children of  $A$ ,  $\text{exp}(\alpha_i)^{\text{rev}}$ , and  $\mathcal{Y}$  contains all the expansions of  
 172 the nonempty rule suffixes,  $\text{exp}(\alpha_{i+1}) \cdots \text{exp}(\alpha_t)$ . Both sets are sorted lexicographically and  
 173 placed on a grid with (less than)  $g$  points,  $t - 1$  for each rule  $A \rightarrow \alpha_1 \cdots \alpha_t$ . Given a pattern  
 174  $P[1..m]$ , for each cut  $P = R | Q$ , we first find the lexicographic ranges  $[s_x, e_x]$  of  $R^{\text{rev}}$  in  $\mathcal{X}$   
 175 and  $[s_y, e_y]$  of  $Q$  in  $\mathcal{Y}$ . Each point  $(x, y) \in [s_x, e_x] \times [s_y, e_y]$  represents a primary occurrence  
 176 of  $P$ . Grid points are augmented with their locus node  $v$  and offset  $|\text{exp}(\alpha_1) \cdots \text{exp}(\alpha_i)|$ .

177 Once we identify the locus node  $v$  (with label  $A$ ) of a primary occurrence, every other  
 178 mention of  $A$  or its ancestors in the grammar tree, and recursively, of the ancestors of those  
 179 mentions, yields a secondary occurrence of  $P$ . Those are efficiently tracked and reported  
 180 [8, 9, 6]. An important *consistency* observation for counting is that the amount of secondary  
 181 occurrences triggered by each primary occurrence is fixed. See Figure 1.

182 The original approach [8, 9] spends time  $O(m^2)$  to find the ranges  $[s_x, e_x]$  and  $[s_y, e_y]$  for  
 183 the  $m - 1$  cuts of  $P$ ; this was later improved to  $O(m \log n)$  [6]. Each primary occurrence

## XX:6 Counting on General Run-Length Grammars

184 found in the grid ranges takes time  $O(\log^\epsilon g)$  using geometric data structures, whereas each  
185 secondary occurrence requires  $O(1)$  time. Overall, the *occ* occurrences of  $P$  in  $T$  are listed in  
186 time  $O(m \log n + \text{occ} \log^\epsilon g)$ .

187 To generalize this solution to RLCFGs [6, App. A.4], rules  $A \rightarrow B^s$  are added as a point  
188  $(x, y) = (\exp(B)^{\text{rev}}, \exp(B)^{s-1})$  in the grid. This suffices to capture every primary occurrence  
189 of the corresponding rule  $A \rightarrow B \cdots B$ : If there are primary occurrences with the cut  
190  $P = R \mid Q$  in  $B \cdots B$ , then one is aligned with the first phrase boundary,  $\exp(B) \mid \exp(B)^{s-1}$ .  
191 Precisely, there is space to place  $Q$  right after the first  $t = s - \lceil |Q|/|B| \rceil$  phrase boundaries.  
192 When the point  $(x, y)$  is retrieved for a given cut, then,  $t$  primary occurrences are declared  
193 with offsets  $|B| - |R|, 2|B| - |R|, \dots, t|B| - |R|$  within  $\exp(A)$ . The amount of secondary  
194 occurrences triggered by each such primary occurrence still depends only on  $A$ .

### 4 Counting with Grammars

196 Navarro [40] obtained the first result in counting the number of occurrences of a pattern  
197  $P[1..m]$  in a text  $T[1..n]$  represented by a CFG of size  $g$ , within time  $O(m^2 + m \log^{2+\epsilon} g)$ , for  
198 any constant  $\epsilon > 0$ , and using  $O(g)$  space. His method relies on the *consistency* observation  
199 above, which allows enhancing the grid described in Section 3 with the number  $c(A)$  of  
200 (primary and) secondary occurrences associated with each point. At query time, for each  
201 pattern cut, one sums the number of occurrences in the corresponding grid range using  
202 the technique mentioned in Section 2.4. The final complexity is obtained by aggregating  
203 over all  $m - 1$  cuts of  $P$  and considering the  $O(m^2)$  time required to identify all the ranges.  
204 Christiansen et al. [6, Thm. A.5] later improved this time to just  $O(m \log n + m \log^{2+\epsilon} g)$ , by  
205 using more modern techniques to find the grid range of all cuts of  $P$ .

206 Christiansen et al. [6] also presented a method to count in  $O(m + \log^{2+\epsilon} n)$  time on a  
207 *particular* RLCFG of size  $g_{rl} = O(\gamma \log(n/\gamma))$ , where  $\gamma$  is the size of the smallest string  
208 attractor [24] of  $T$ . They also show that by increasing the space to  $O(\gamma \log(n/\gamma) \log^\epsilon n)$  one  
209 can reach the optimal counting time,  $O(m)$ . The grammar properties allow reducing the  
210 number of cuts of  $P$  to check to  $O(\log m)$ , instead of the  $m - 1$  cuts used on general RLCFGs.

211 Christiansen et al. build on the same idea of enhancing the grid with the number of  
212 secondary occurrences, but the process is considerably more complex on RLCFGs, because  
213 the consistency property exploited by Navarro [40] does not hold on run-length rules  $A \rightarrow B^s$ :  
214 the number of occurrences triggered by a primary occurrence with cut  $P = R \mid Q$  found from  
215 the point  $(\exp(B)^{\text{rev}}, \exp(B)^{s-1})$  depends on  $s, |B|$ , and  $|R|$ . Their counting approach relies  
216 on another property that is specific of their RLCFG [6, Lem. 7.2]:

217 ► **Property 1.** *For every run-length rule  $A \rightarrow B^s$ , the shortest period of  $\exp(A)$  is  $|B|$ .*

218 This property facilitates the division of the counting process into two cases. For each  
219 run-length rule  $A \rightarrow B^s$ , they introduce two points,  $(x, y') = (\exp(B)^{\text{rev}}, \exp(B))$  and  
220  $(x, y'') = (\exp(B)^{\text{rev}}, \exp(B)^2)$ , in the grid. These points are associated with the values  $c(A)$   
221 and  $(s - 2) \cdot c(A)$ , respectively. The counting process is as follows: for a cut  $P = R \mid Q$ , if  
222  $Q \sqsubseteq \exp(B)$ , then it will be counted  $c(A) + (s - 2) \cdot c(A) = (s - 1) \cdot c(A)$  times, as both points  
223 will be within the search range. If  $Q$  instead exceeds  $\exp(B)$ , but still  $Q \sqsubseteq \exp(B)^2$ , then it  
224 will be counted  $(s - 2) \cdot c(A)$  times, solely by point  $(x, y'')$ . Finally if  $Q$  exceeds  $\exp(B)^2$ ,  
225 then  $Q$  is periodic (with  $p(Q) = |B|$ ).

226 They handle that remaining case as follows. Given a cut  $P = R \mid Q$  and the period  
227  $p = p(Q) = |B|$ , where  $|Q| > 2p$ , the number of primary occurrences of this cut inside rule  
228  $A \rightarrow B^s$  is  $s - \lceil |Q|/p \rceil$  (cf. the end of Section 3). Let  $D$  be the set of rules  $A \rightarrow B^s$  within

229 the grid range of the cut, and  $c(A)$  the number of (primary and secondary) occurrences of  $A$ .  
 230 Then, the number of occurrences triggered by the primary occurrences found within symbols  
 231 in  $D$  for this cut is

$$232 \quad \sum_{A \rightarrow B^s \in D} c(A) \cdot s - c(A) \cdot \lceil |Q|/p \rceil.$$

233 For each run-length rule  $A \rightarrow B^s$ , they compute a Karp–Rabin signature (Section 2.3)  
 234  $\kappa(\exp(B))$  and store it in a perfect hash table, associated with values

$$235 \quad C(B, s) = \sum \{c(A) : A \rightarrow B^s, s' \geq s\},$$

$$236 \quad C'(B, s) = \sum \{s' \cdot c(A) : A \rightarrow B^{s'}, s' \geq s\}.$$

237 Additionally, for each such  $B$ , the authors store the set  $s(B) = \{s : A \rightarrow B^s\}$ .

238 At query time, they calculate the shortest period  $p = p(P)$ . For each cut  $P = R | Q$ ,  
 239  $Q$  is periodic if  $|Q| > 2p$ . If so, they compute  $k = \kappa(Q[1..p])$ , and if there is an entry  $B$   
 240 associated with  $k$  in the hash table, they add to the number of occurrences found up to then

$$241 \quad C'(B, s_{min}) - C(B, s_{min}) \cdot \lceil |Q|/p \rceil,$$

242 where  $s_{min} = \min\{s \in s(B), (s-1) \cdot |B| \geq |Q|\}$  is computed using exponential search over  
 243  $s(B)$  in  $O(\log m)$  time. Note that they exploit the fact that *the number of repetitions to*  
 244 *subtract*,  $\lceil |Q|/p \rceil$ , depends only on  $p = |B|$ , and not on the exponent  $s$  of rules  $A \rightarrow B^s$ .

245 The total counting time, on a grammar of size  $g_{rl}$ , is  $O(m \log n + m \log^{2+\epsilon} g_{rl})$ . In their  
 246 particular grammar, the number of cuts to consider is  $O(\log m)$ , which allows reducing the  
 247 cost of computing the grid ranges to  $O(m)$ . The signatures of all substrings of  $P$  are also  
 248 computed in  $O(m)$  time, as mentioned in Section 2.3. Considering the grid searches, the total  
 249 cost for counting the pattern occurrences drops to  $O(m + \log^{2+\epsilon} g_{rl}) \subseteq O(m + \log^{2+\epsilon} n)$ .

250 Recently, Kociumaka et al. [27] employed this same approach to count the occurrences  
 251 of a pattern in a smaller RLCFG that uses  $O(\delta \log \frac{n \log |\Sigma|}{\delta \log n})$  space, where  $\delta \leq \gamma$ . They  
 252 demonstrated that the RLCFG they produce satisfies Property 1 [6, Lem. 7.2], which is  
 253 necessary to apply the described scheme.

## 254 5 Our Solution

255 We now describe a solution to count the occurrences in arbitrary RLCFGs, where the  
 256 convenient Property 1 used in the literature may not hold. We start with a simple observation.

257 ► **Lemma 3.** *Let  $A \rightarrow B^s$  be a rule in a RLCFG. Then  $p(A)$  divides  $|B|$ .*

258 **Proof.** Clearly  $|B|$  is a period of  $\exp(A)$  because  $\exp(A) = \exp(B)^s$ . By Lemma 2, then,  
 259 since  $|B| \leq |A|/2$ ,  $p(A)$  divides  $|B|$ . ◀

260 Some parts of our solution make use of the shortest period of  $\exp(A)$ . We now define  
 261 some related notation.

262 ► **Definition 4.** *Given a rule  $A \rightarrow B^s$  with  $s \geq 2$ , let  $p = p(A)$  (which divides  $|B|$  by Lemma  
 263 3). The corresponding transformed rule is  $A \rightarrow \hat{B}^{s'}$ , where  $\hat{B}$  is a new nonterminal such that  
 264  $\exp(\hat{B}) = \exp(A)[1..p]$ , and  $s' = s \cdot (|B|/p)$ .*

265 There seems to be no way to just transform all run-length rules (which would satisfy  
 266 Property 1,  $p(A) = |\hat{B}|$ ) without blowing up the RLCFG size by a logarithmic factor. We  
 267 will use another approach instead. We classify the rules into two categories.

## XX:8 Counting on General Run-Length Grammars

268 ► **Definition 5.** Given a rule  $A \rightarrow B^s$  with  $s \geq 2$ , we say that  $A$  is of type-E (for Equal) if  
269  $p(A) = |\hat{B}| = |B|$ ; otherwise,  $p(A) = |\hat{B}| < |B|$  and we say that  $A$  is of type-L (for Less).

270 We build on Navarro's solution [40] for counting on CFGs, which uses an enhanced grid  
271 where points count all the occurrences they trigger. The grid ranges are found with the more  
272 recent technique [6] that takes  $O(m \log n)$  time. Further, we treat type-E rules exactly as  
273 Christiansen et al. [6] handle the run-length rules in their specific RLCFGs, as described  
274 in Section 4. This is possible because type-E rules, by definition, satisfy Property 1. Their  
275 method, however, assumes that no two symbols  $B \neq B'$  have the same expansion. To relax  
276 this assumption, symbols  $B$  with the same expansion should collectively contribute to the  
277 same entries of  $C(\cdot, s)$  and  $C'(\cdot, s)$ . We thus index those tables using  $\kappa(\exp(B))$  rather than  
278  $B$ , and for simplicity write  $C(\pi, s)$ ,  $C'(\pi, s)$ , and  $s(\pi)$ , where  $\pi = \exp(B)$ .

279 Since each primary occurrence is found in exactly one rule, we can decompose the process  
280 of counting by adding up the occurrences found inside type-E and type-L rules. We are then  
281 left with the more complicated problem of counting occurrences found from type-L rules.  
282 We start with another observation.

283 ► **Observation 6.** If  $A \rightarrow B^s$  is a type-L rule, then  $|B| \geq 2|\hat{B}|$

284 **Proof.** If  $A$  is a type-L rule then  $p(A) = |\hat{B}| < |B|$ . In addition, by Lemma 3,  $|\hat{B}|$  divides  
285  $|B|$ . Therefore  $|B| \geq 2|\hat{B}|$  ◀

286 For type-L rules, we will generalize the strategy of Section 4: the cases where  $|Q| \leq 2|\hat{B}|$   
287 will be handled by adding points to the enhanced grid; in the other cases we will use new data  
288 structures that exploit the fact (to be proved) that  $Q$  is periodic. Note that each partition  
289  $P = R \mid Q$  may correspond to different cases for different run-length rules, so our technique  
290 will consider all the cases for each partition.

### 291 5.1 Case $|Q| \leq 2|\hat{B}|$

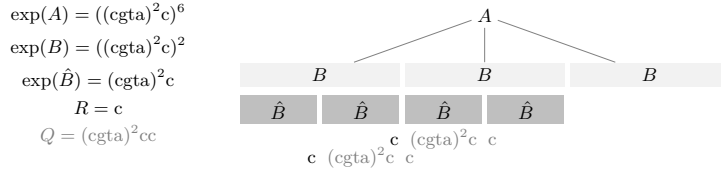
292 To capture the primary occurrences with cut  $P = R \mid Q$  inside type-L rules  $A \rightarrow B^s$  where  
293  $|Q| \leq 2|\hat{B}|$ , we will incorporate the points  $(x_p, y'_p) = (\exp(\hat{B})^{\text{rev}}, \exp(\hat{B}))$  and  $(x_p, y''_p) =$   
294  $(\exp(\hat{B})^{\text{rev}}, \exp(\hat{B})^2)$  into the enhanced grid outlined in Sections 3 and 4, assigning the values  
295  $-(s-1) \cdot c(A)$  and  $2 \cdot (s-1) \cdot c(A)$  to each, respectively. The point  $(x_p, y'_p)$  will capture  
296 the occurrences where  $|R|, |Q| \leq |\hat{B}|$ . Note that these occurrences will also find the point  
297  $(x_p, y''_p)$ , so the final result will be  $(2-1) \cdot (s-1) \cdot c(A) = (s-1) \cdot c(A)$ .

298 The point  $(x_p, y''_p)$  will also account for the primary occurrences where  $|R| \leq |\hat{B}|$  and  
299  $|\hat{B}| < |Q| \leq 2|\hat{B}|$ . Observation 6 establishes that  $|B| \geq 2|\hat{B}|$ , so for each such primary  
300 occurrence of cut  $R \mid Q$ , with offset  $j$  in  $\exp(A)$ , there is a second primary occurrence at  
301  $j - |\hat{B}|$  with cut  $P = R' \mid Q'$ , where  $|\hat{B}| < |R'| = |R| + |\hat{B}| \leq 2|\hat{B}|$  and  $|Q'| = |Q| - |\hat{B}| \leq |\hat{B}|$ .  
302 This second cut will not be captured by the points we have inserted because  $|R'| > |\hat{B}|$ . The  
303 other occurrences where  $P$  matches to the left of  $j - |\hat{B}|$  fall within  $B$  (and thus are not  
304 primary), because we already have  $|Q'| \leq |\hat{B}|$  in this second occurrence. Thus, for each of  
305 the  $s$  copies of  $B$  (save the last), we will have two primary occurrences. This yields a total of  
306  $2 \cdot (s-1) \cdot c(A)$  occurrences, which are properly counted in the points  $(x_p, y''_p)$ . See Figure 2.

### 307 5.2 Case $|Q| > 2|\hat{B}|$

308 We first show that, for  $Q$  to be longer than  $2|\hat{B}|$  in some run-length rule,  $P$  must be periodic.





■ **Figure 2** We show the occurrences captured by the point  $(x_p, y_p'') = (\exp(\hat{B}), \exp(\hat{B})^2)$ . Note how the occurrence in the first row is correctly captured by  $(x_p, y_p'')$ , whereas that in the second row is not captured by any point. Consequently, the first row is effectively counted twice. Given that the point  $(x_p, y_p'')$  is assigned a weight of  $2 \cdot (s-1) \cdot c(A)$ , the total number of occurrences is  $4 \cdot c(A)$ .

309 ► **Lemma 7.** Let  $P$ , with  $p = p(P)$ , have a primary occurrence with cut  $P = R \mid Q$  in the  
310 rule  $A \rightarrow B^s$ , with  $p(A) = |\hat{B}|$  and  $|Q| > 2|\hat{B}|$ . Then it holds that  $p = p(A)$ .

311 **Proof.** Since  $|P| \geq |\hat{B}|$  and  $P$  is contained within  $\exp(A) = \exp(\hat{B})^{s'}$ , by branch 3 of  
312 Definition 1,  $|\hat{B}|$  must be a period of  $P$ . Thus,  $p = p(P) \leq |\hat{B}|$ . Suppose, for contradiction,  
313 that  $p < |\hat{B}|$ . According to Lemma 2, because  $|\hat{B}| \leq |Q|/2 \leq |P|/2$  is a period of  $P$ , it  
314 follows that  $p$  divides  $|\hat{B}|$ . Since  $\exp(\hat{B})$  is contained in  $P$ , again by branch 3 of Definition 1  
315 it follows that  $p < |\hat{B}| \leq |B|$  is a period of  $\exp(B)$ , and thus of  $\exp(A)$ , contradicting the  
316 assumption that  $p(A) = |\hat{B}|$ . Hence, we conclude that  $p = |\hat{B}|$ . ◀

317 Note that  $P$  is then periodic because  $p(P) = p(A) = |\hat{B}| < |Q|/2 \leq |P|/2$ , and  $Q$  is  
318 also periodic by branch 3 of Def. 1, because it occurs inside  $P$  and  $|Q| \geq 2p$ . The following  
319 definition will help us show that we capture every primary occurrence exactly once.

320 ► **Definition 8.** The alignment of a primary occurrence  $x$  found with cut  $P = R \mid Q$  inside  
321 the type- $L$  rule  $A \rightarrow B^s$  is  $\text{align}(x) = 1 + ((|R| - 1) \bmod |\hat{B}|)$ .

322 The definition is sound because every primary occurrence is found using exactly one  
323 cut  $P = R \mid Q$ . Note that  $\text{align} \in [1..|\hat{B}|]$  is the distance from the starting position of  
324 an occurrence, within  $\exp(A)$ , to the start of the next copy of  $\exp(\hat{B})$ . We will explore  
325 all the possible cuts of  $P$ , but each rule  $A \rightarrow B^s$  will be probed only with the cuts where  
326  $1 \leq |R| \leq |\hat{B}|$ . From those cuts, all the corresponding primary occurrences aligned with the  
327  $s' - 1$  boundaries between copies of  $\hat{B}$  (i.e., with the same alignment,  $|R|$ ) will be captured.

328 We distinguish two subcases, depending on whether  $Q$  is longer than  $B$  or not. If it is,  
329 we must ensure that in the alignments we count the occurrence is fully within  $\exp(A)$ . If it  
330 is not, we must ensure that the alignments we count do correspond to primary occurrences  
331 (i.e., they cross a border between copies of  $B$ ).

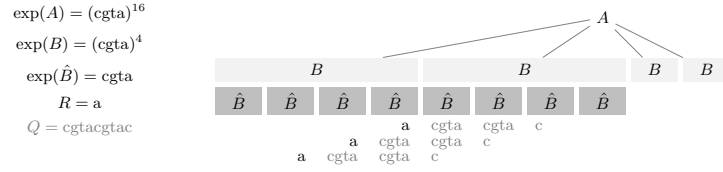
### 332 5.2.1 Case $2|\hat{B}| < |Q| \leq |B|$

333 To handle this case, we construct a specific data structure based on the period  $|\hat{B}|$ . The  
334 proposed solution is supported by the following lemma.

335 ► **Lemma 9.** Let  $P$ , with  $p = p(P)$ , have a primary occurrence with cut  $P = R \mid Q$  in the  
336 type- $L$  rule  $A \rightarrow B^s$ , with  $p(A) = |\hat{B}|$ ,  $|R| \leq |\hat{B}|$ , and  $2|\hat{B}| < |Q| \leq |B|$ . Then, the number  
337 of primary occurrences of  $P$  in  $\exp(A)$  is  $(s-1) \cdot \lceil |Q|/p \rceil$ .

338 **Proof.** Since  $|R| \leq |\hat{B}|$ ,  $R$  can be aligned at the end of the  $|B|/|\hat{B}|$  positions where  $\exp(\hat{B})$   
339 starts in  $\exp(B)$ . No other alignments are possible for the cut  $R \mid Q$  because, by Lemma 7,  
340  $p = |\hat{B}|$  and another alignment would imply that  $P$  aligns with itself with an offset smaller  
341 than  $p$ , a contradiction by branch 2 of Definition 1.

## XX:10 Counting on General Run-Length Grammars



■ **Figure 3** If  $2|\hat{B}| < |Q| \leq |B|$ , there are  $\lceil |Q|/p \rceil$  primary occurrences around the boundary between any two blocks  $B$  (we zoom on one) with the cut  $P = R | Q$ . We show the possible alignments of  $P$  below the blocks  $\hat{B}$ . For a rule  $A \rightarrow B^s$  there are  $(s - 1)$  boundaries, yielding  $(s - 1) \cdot \lceil |Q|/p \rceil$  primary occurrences. In this case,  $\lceil |Q|/p \rceil = 3$  and  $s - 1 = 3$ , yielding 9 primary occurrences.

342 Those alignments correspond to primary occurrences only if  $P$  does not fall completely  
 343 within  $\exp(B)$ . The alignments that correspond to primary occurrences are then those where  
 344  $R$  is aligned at the end of the last  $\lceil |Q|/|\hat{B}| \rceil$  ending positions of copies of  $\hat{B}$ , all of which  
 345 start within  $\exp(B)$  because  $|Q| \leq |B|$ . This is equivalent to  $\lceil |Q|/p \rceil$ , as  $p = |\hat{B}|$  by Lemma  
 346 7. Thus, the number of primary occurrences of  $P$  in  $A$  is  $(s - 1) \cdot \lceil |Q|/p \rceil$ . See Figure 3. ◀

347 Based on Lemma 9 we introduce our first period-based data structure. Considering the  
 348 solution described in Section 4, where Property 1 holds, the challenge with type-L rules  
 349  $A \rightarrow B^s$  (i.e., rules that differ from their transformed version  $A \rightarrow \hat{B}^{s'}$ ) is that the number of  
 350 alignments with cut  $R | Q$  inside  $\exp(A)$  is  $s' - \lceil |Q|/p \rceil$ , but  $B$  does not determine  $p = p(A)$ .  
 351 We will instead use  $\hat{B}$  to index those nonterminals  $A$ .

352 For each type-L rule  $A \rightarrow B^s$  ( $A \rightarrow \hat{B}^{s'}$  being its transformed version), we compute its  
 353 signature  $\kappa(\exp(\hat{B}))$  (recall Section 2.3) and store it in a perfect hash table  $H$ . Each entry in  
 354 table  $H$ , which corresponds to a specific signature  $\kappa(\pi)$ , will be linked to an array  $F_\pi$ . Each  
 355 position  $F_\pi[i]$  represents a type-L rule  $A_i \rightarrow B_i^{s_i}$  where  $\kappa(\exp(\hat{B}_i)) = \kappa(\pi)$ . The rules  $A_i$  are  
 356 sorted in  $F_\pi$  by decreasing lengths  $|B_i|$ . We also store a field with the cumulative sum

$$357 \quad F_\pi[i].sum = \sum_{1 \leq j \leq i} (s_j - 1) \cdot c(A_j).$$

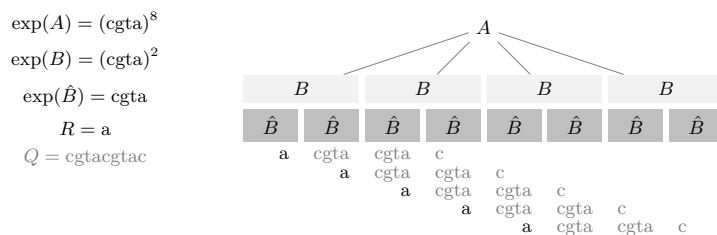
358 Given a pattern  $P[1..m]$ , we first calculate its shortest period  $p = p(P)$ . For each cut  
 359  $P = R | Q$  with  $1 \leq |R| \leq \min(p, m - 2p - 1)$ , we compute  $\kappa(\pi)$  for  $\pi = Q[1..p]$  to identify  
 360 the corresponding array  $F_\pi$  in  $H$ . Note that we only consider the cuts  $R | Q$  where  $|R| \leq p$ ,  
 361 as this corresponds precisely to  $|R| \leq |\hat{B}|$  for the rules stored in  $F_\pi$ ; note  $p = |\pi|$ . In addition,  
 362 the condition  $|R| \leq m - 2p - 1$  ensures that  $|Q| > 2p = 2|\hat{B}|$ , so we are correctly enforcing the  
 363 condition of this subsection and focusing on the occurrences of each alignment  $align = |R|$   
 364 one by one. We will find in  $H$  every (transformed) rule  $A \rightarrow \hat{B}^{s'}$  where  $\hat{B} = \pi$ , sharing the  
 365 period  $p$  with  $Q$ , as well as its prefix  $\pi = \exp(B)[1..p] = Q[1..p]$ . Once we have obtained  
 366 the array  $F_\pi$ , we find the largest  $i$  such that  $|B_i| \geq |Q|$ . The number of primary occurrences  
 367 for the cut  $P = R | Q$  in type-L rules where  $2|\hat{B}| < |Q| \leq |B|$  is then  $F_\pi[i].sum \cdot \lceil |Q|/p \rceil$ .

### 368 5.2.2 Case $|Q| > |B|$

369 Our analysis for the remaining case is grounded on the following lemma.

370 ► **Lemma 10.** *Let  $P$ , with  $p = p(P)$ , have a primary occurrence in a type-L rule  $A \rightarrow B^s$   
 371 with cut  $P = R | Q$ , with  $|R| \leq p$  and  $|Q| > |B|$ . Then it holds that  $p = p(A)$  and  $|Q| > 2p$ .*

372 **Proof.** If  $A$  is a type-L rule and  $P$  has an occurrence within  $A$  such that  $|Q| > |B|$ , then  
 373 we have  $|Q| > |B| \geq 2|\hat{B}|$  (by Observation 6). Since we can express  $A$  as  $A \rightarrow \hat{B}^{s'}$ , we can  
 374 similarly use Lemma 7 to conclude that  $p = p(A) = |\hat{B}|$ ; further,  $|Q| > 2p$ . ◀



■ **Figure 4** If  $|Q| > |B|$ , we can compute all occurrences of  $P$  around blocks  $\hat{B}$  without the risk of any occurrence being fully contained in a block  $B$ : the number of primary occurrences of  $P$  in  $\exp(A)$  is simply  $s' - \lceil |Q|/p \rceil$ . In this example, with  $s' = 8$  and  $\lceil |Q|/p \rceil = 3$ , there are 5 occurrences.

375 Analogously to Lemma 7, Lemma 10 establishes that, when  $Q$  is sufficiently long, it holds  
 376 that  $p(P) = p(A)$ , so all pertinent rules of the form  $A \rightarrow B^s$  can be classified according to  
 377 their minimal period,  $p(A)$ . This period coincides with  $p = p(P)$  when  $P$  has an occurrence  
 378 in a type-L rule such that  $|Q| > |B|$ . Further,  $|Q| > 2p$ .

379 We also need an analogous to Lemma 9 for the case  $|Q| > |B|$ ; this is given next.

380 ► **Lemma 11.** *Let  $P$ , with  $p = p(P)$ , have a primary occurrence with cut  $P = R | Q$  in*  
 381 *the type-L rule  $A \rightarrow B^s$ , with  $p(A) = |\hat{B}|$ ,  $|R| \leq |\hat{B}|$ , and  $|Q| > |B|$ . Then, the number of*  
 382 *primary occurrences of  $P$  in  $\exp(A)$  is  $s' - \lceil |Q|/p \rceil$ .*

383 **Proof.** Since  $|R| \leq |\hat{B}|$ ,  $R$  can be aligned at the end of the  $s'$  positions where  $\exp(\hat{B})$  starts in  
 384  $\exp(A)$ . By the same argument of the proof of Lemma 9, no other alignments are possible for  
 385 the cut  $R | Q$ . Unlike in Lemma 9, all those alignments correspond to primary occurrences,  
 386 because  $Q$  is always long enough to exceed  $B$ . Also unlike in Lemma 9,  $Q$  may exceed  $A$ ,  
 387 in which case the occurrence must not be counted in this rule. The alignments that must  
 388 not be counted are then those where  $R$  is aligned at the end of the last  $\lceil |Q|/|\hat{B}| \rceil$  ending  
 389 positions of copies of  $\hat{B}$ . This is equivalent to  $\lceil |Q|/p \rceil$ , as  $p = |\hat{B}|$  by Lemma 10. Thus, the  
 390 number of primary occurrences of  $P$  in  $A$  is  $s' - \lceil |Q|/p \rceil$ . See Figure 4. ◀

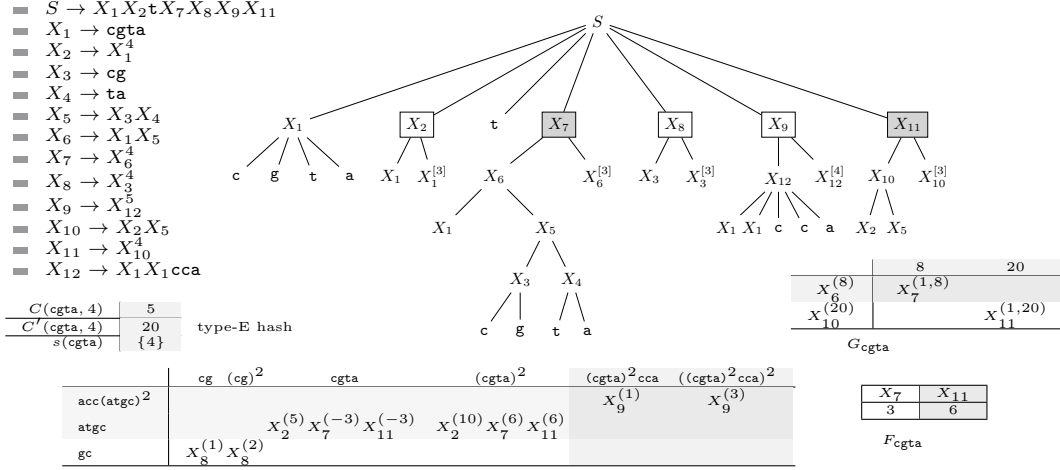
391 We then enhance table  $H$ , introduced in Section 5.2.1, with a second period-based data  
 392 structure. Each entry in table  $H$ , corresponding to some  $\kappa(\pi)$ , will additionally store a grid  
 393  $G_\pi$ . In this grid, each row represents a type-L rule  $A \rightarrow B^s$  whose transformed version is  
 394  $A \rightarrow \hat{B}^{s'}$ , that is, such that  $\pi = \exp(\hat{B}) = \exp(B)[1..p]$ . The rows are sorted by increasing  
 395 lengths  $|B|$  (note  $|B| \geq |\pi| = p$  for all  $B$  in  $G_\pi$ ). The columns represent the different  
 396 exponents  $s'$  of the transformed rules. The row of rule  $A \rightarrow B^s$  has then a unique point at  
 397 column  $s'$ , and we associate two values with it:  $C'_\pi(A) = c(A)$  and  $C''_\pi(A) = s' \cdot c(A)$ . Since  
 398 no rule appears in more than one grid, the total space for all grids is in  $O(g_{rl})$ .

399 Given a pattern  $P[1..m]$ , we proceed analogously as explained at the end of Section 5.2.1  
 400 in order to identify  $F_\pi$ : We compute  $p = p(P)$ , and for each cut  $P = R | Q$  with  $1 \leq |R| \leq$   
 401  $\min(p, m - 2p - 1)$ , we calculate  $\kappa(\pi)$ , for  $\pi = Q[1..p]$ , to find the corresponding grid  $G_\pi$  in  
 402  $H$ . On the type-L rules  $A \rightarrow B^s$ , this tries out every possible alignment  $align = |R|$ , one by  
 403 one, from 1 to  $|\hat{B}|$ . The limit  $|R| < m - 2p$  can also be set because, by Lemma 10, it must  
 404 hold  $|Q| > 2|\hat{B}|$  on the rules of  $G_\pi$  we find with the cut  $P = R | Q$ .

405 We must enforce two conditions on the rules of  $G_\pi$  to consider: (a)  $|Q| > |B|$  as  
 406 corresponds in this subsection, and (b)  $s' - \lceil |Q|/p \rceil \geq 0$ , that is,  $Q$  fits within  $\exp(A)$ . The  
 407 complying rules then contribute  $c(A) \cdot (s' - \lceil |Q|/p \rceil) = C''_\pi(A) - C'_\pi(A) \cdot \lceil |Q|/p \rceil$  by Lemma 11.

408 To enforce those conditions, we find in  $G_\pi$  the largest row  $y$  representing a rule  $A \rightarrow B^s$   
 409 such that  $|B| < |Q|$ . We also find the smallest column  $x$  where  $(s' =) x \geq \lceil |Q|/p \rceil$ . The

## XX:12 Counting on General Run-Length Grammars



■ **Figure 5** On top, a RLCFG on the left and its grammar tree on the right. Type-E rules are enclosed in white rectangles and Type-L rules in gray rectangles. Below the rules we show the values  $C(B, s)$  and  $C'(B, s)$  [6] we use to handle the E-type rules (see Section 4); we only show those for  $\text{exp}(X_1) = \text{cgta}$ . On the bottom left we show the points we add to the standard grid. The points for type-E rules are represented as  $A^{(c(A))}$  and  $A^{((s-2) \cdot c(A))}$  and those for type-L rules as  $A^{(-(s-1) \cdot c(A))}$  and  $A^{(2 \cdot (s-1) \cdot c(A))}$ . The bottom right shows the grid  $G_\pi$  and the array  $F_\pi$  for the transformed rules  $A \rightarrow \hat{B}^{s'}$  where  $\hat{B} = \pi = \text{cgta}$ . In  $F_\pi$  we show the fields  $F[i].\text{sum}$ . In  $G_\pi$ , the row labels show  $B^{(|B|)}$  and the column labels show  $s'$ ; the points show  $A^{(C', C')}$ . Consider the cut  $P = \text{a} \mid \text{cgta} \text{c} \text{gta} \text{c}$ , with  $p(P) = 4$ . We identify 9 occurrences in type-E rules: 4 are found within the rule  $X_9$  using the standard grid, while the remaining 5 are determined via the values of  $C(X_1, s)$  and  $C'(X_1, s)$ . These 5 occurrences specifically arise within  $\text{exp}(X_2) = (\text{cgta})^4$ . Similarly, in the type-L rules, we detect 14 occurrences: 9 occur within the rule  $X_{11}$ , identified using the  $F_{\text{cgta}}$  array, and the remaining 5 arise within  $\text{exp}(X_7) = (\text{cgta})^8$ , captured using the  $G_{\text{cgta}}$  grid. The final two occurrences of this cut are located using standard CFG rules at  $\text{exp}(S)[4..13]$  ( $X_1 \cdot X_2$ ) and  $\text{exp}(S)[111..120]$  ( $X_9 \cdot X_{11}$ ).

410 points in the range  $[x, n] \times [1, y]$  of the grid then correspond to the set  $D$  of type-L run-length  
 411 rules where we have a primary occurrence with  $|Q| > |B|$ . We aggregate the values  $C'_\pi$  and  
 412  $C''_\pi$  from the range, which yields the correct sum of all the pertinent occurrences:

$$413 \quad \left( \sum_{A \rightarrow B^s \in D} C''_\pi(A) \right) - \left( \sum_{A \rightarrow B^s \in D} C'_\pi(A) \right) \cdot \lceil |Q|/p \rceil = \sum_{A \rightarrow B^s \in D} c(A) \cdot s' - c(A) \cdot \lceil |Q|/p \rceil.$$

414 Figure 5 gives a thorough example.

### 415 5.3 The final result

416 Our structure extends the grid of Section 4, built for non-run-length rules, with one point per  
 417 run-length rule: those of type-E are handled as described in Section 4 and those of type-L as  
 418 in Section 5. Thus the structure is of size  $O(g_{rl})$  and range queries on the grid take time  
 419  $O(\log^{2+\epsilon} g_{rl})$ . Occurrences on such a grid are counted in time  $O(m \log n + m \log^{2+\epsilon} g_{rl})$  [6,  
 420 Thm. A.5]. This is also the time to count the occurrences in type-E rules for our solution,  
 421 and those in type-L rules when  $|Q| \leq 2|B_p|$  (Section 5.1).

422 For our period-based data structures (Sections 5.2.1 and 5.2.2), we calculate  $p(P)$  in  
 423  $O(m)$  time [10], and compute all prefix signatures of  $P$  in  $O(m)$  time as well, so that later  
 424 any substring signature is computed in  $O(1)$  time (Section 2.3). The limits in the arrays  $F_\pi$   
 425 and in the grids  $G_\pi$  can be found with exponential search in time  $O(\log m)$  (we might need

426 to group rows/columns with identical values to achieve this time). The range sums for  $C'_\pi$   
 427 and  $C''_\pi$  take time  $O(\log^{2+\epsilon} g_{rl})$ . They are repeated for each of the  $O(m)$  cuts of  $P$ , adding  
 428 up to time  $O(m \log^{2+\epsilon} g_{rl})$ . Those are then within the previous time complexities as well.

429 ► **Theorem 12.** *Let a RLCFG of size  $g_{rl}$  represent a text  $T[1..n]$ . Then, for any constant*  
 430  *$\epsilon > 0$ , we can build in  $O(n \log n)$  expected time an index of size  $O(g_{rl})$  that counts the number*  
 431 *of occurrences of a pattern  $P[1..m]$  in  $T$  in time  $O(m \log n + m \log^{2+\epsilon} g_{rl}) \subseteq O(m \log^{2+\epsilon} n)$ .*

432 Just as for previous schemes [6], the construction time is dominated by the  $O(n \log n)$   
 433 expected time to build the collision-free Karp–Rabin functions [3].

434 Note that the bulk of the search cost are the geometric queries, which are easily done in  
 435  $O(\log n)$  time if we store cumulative sums in all the levels of the data structure [5, 40]. More  
 436 generally, setting Navarro’s  $\epsilon$  to  $1/\log^{1-\delta} g_{rl}$  [40, Thm. 3], we obtain the following tradeoff.

437 ► **Corollary 13.** *Let a RLCFG of size  $g_{rl}$  represent a text  $T[1..n]$ . Then, for any constant  $0 \leq$   
 438  $\delta < 1$ , we can build in  $O(n \log n)$  expected time an index of size  $O(g_{rl} \log^{1-\delta} g_{rl})$  that counts*  
 439 *the occurrences of a pattern  $P[1..m]$  in  $T$  in time  $O(m \log n + m \log^{1+\delta} g_{rl}) \subseteq O(m \log^{1+\delta} n)$ .*

## 440 5.4 An application

441 Recent work [18, 37] shows how to compute the maximal exact matches (MEMs) of  $P[1..m]$   
 442 in  $T[1..n]$ , which are the maximal substrings of  $P$  that occur in  $T$ , in case  $T$  is represented  
 443 with an arbitrary RLCFG. Navarro [41] extends the results to  $k$ -MEMs, which are maximal  
 444 substrings of  $P$  that occur at least  $k$  times in  $T$ . To obtain good time complexities for large  
 445 enough  $k$ , he resorts to counting occurrences of substrings  $P[i..j]$  with the grammar. His  
 446 Thm. 7, however, works only for CFGs, as no efficient counting algorithm existed on RLCFGs.  
 447 In turn, his Thm. 8 works only for a particular RLCFG. We can now state his result on an  
 448 arbitrary RLCFG; by his Thm. 11 this also extends to “ $k$ -rare MEMs”.

449 ► **Corollary 14** (cf. [41, Thm. 7]). *Let a RLCFG of size  $g_{rl}$  generate only  $T[1..n]$ . Then,*  
 450 *for any constant  $\epsilon > 0$ , we can build a data structure of size  $O(g_{rl})$  that finds the  $k$ -MEMs*  
 451 *of any given pattern  $P[1..m]$ , for any  $k > 0$  given with  $P$ , in time  $O(m^2 \log^{2+\epsilon} g_{rl})$ .*

## 452 6 Conclusion

453 We have presented the first solution to the problem of counting the occurrences of a pattern  
 454 in a text represented by an arbitrary RLCFG, which was posed by Christiansen et al. [6]  
 455 in 2020 and solved only for particular cases. This required combining solutions to CFGs  
 456 [40] and particular RLCFGs [6], but also new insights for the general case. The particular  
 457 existing solutions required that  $|B|$  is the shortest period of  $\exp(A)$  in rules  $A \rightarrow B^s$ . While  
 458 this does not hold in general RLCFGs, we proved that, except in some borderline cases  
 459 that can be handled separately, the shortest periods of the pattern and of  $\exp(A)$  must  
 460 coincide. While the particular solutions could associate  $\exp(B)$  with the period of the pattern,  
 461 we must associate many strings  $\exp(A)$  that share the same shortest period, and require  
 462 a more sophisticated geometric data structure to collect only those that qualify for our  
 463 search. Despite those complications, however, we manage to define a data structure of size  
 464  $O(g_{rl})$  from a RLCFG of size  $g_{rl}$ , that counts the occurrences of  $P[1..m]$  in  $T[1..n]$  in time  
 465  $O(m \log^{2+\epsilon} n)$  for any constant  $\epsilon > 0$ , the same result that existed for the simpler case of  
 466 CFGs. Our approach extends the applicability of arbitrary RLCFGs to cases where only  
 467 CFGs could be used, equalizing the available tools to handle both types of grammars.

## 468 — References

- 469 1 Bille, P., Ettiienne, M.B., Gørtz, I.L., Vildhøj, H.W.: Time-space trade-offs for Lempel-Ziv  
470 compressed indexing. *Theoretical Computer Science* **713**, 66–77 (2018)
- 471 2 Bille, P., Landau, G.M., Raman, R., Sadakane, K., Rao, S.S., Weimann, O.: Random access  
472 to grammar-compressed strings and trees. *SIAM Journal on Computing* **44**(3), 513–539 (2015)
- 473 3 Bille, P., Gørtz, I.L., Sach, B., Vildhøj, H.W.: Time-space trade-offs for longest common  
474 extensions. *Journal of Discrete Algorithms* **25**, 42–50 (2014)
- 475 4 Charikar, M., Lehman, E., Liu, D., Panigrahy, R., Prabhakaran, M., Sahai, A., Shelat, A.:  
476 The smallest grammar problem. *IEEE Transactions on Information Theory* **51**(7), 2554–2576  
477 (2005)
- 478 5 Chazelle, B.: A functional approach to data structures and its use in multidimensional  
479 searching. *SIAM Journal on Computing* **17**(3), 427–462 (1988)
- 480 6 Christiansen, A.R., Ettiienne, M.B., Kociumaka, T., Navarro, G., Prezza, N.: Optimal-time  
481 dictionary-compressed indexes. *ACM Transactions on Algorithms (TALG)* **17**(1), 1–39 (2020)
- 482 7 Claude, F., Navarro, G.: Self-indexed grammar-based compression. *Fundamenta Informaticae*  
483 **111**(3), 313–337 (2010)
- 484 8 Claude, F., Navarro, G.: Improved grammar-based compressed indexes. In: *Proc. 19th*  
485 *International Symposium on String Processing and Information Retrieval (SPIRE)*. pp. 180–  
486 192 (2012)
- 487 9 Claude, F., Navarro, G., Pacheco, A.: Grammar-compressed indexes with logarithmic search  
488 time. *Journal of Computer and System Sciences* **118**, 53–74 (2021)
- 489 10 Crochemore, M., Rytter, W.: *Jewels of stringology: text algorithms*. World Scientific (2002)
- 490 11 Ferrada, H., Gagie, T., Hirvola, T., Puglisi, S.J.: Hybrid indexes for repetitive datasets.  
491 *Philosophical Transactions of the Royal Society A* **372**(2016), article 20130137 (2014)
- 492 12 Ferrada, H., Kempa, D., Puglisi, S.J.: Hybrid indexing revisited. In: *Proc. 20th Workshop on*  
493 *Algorithm Engineering and Experiments (ALENEX)*. pp. 1–8 (2018)
- 494 13 Fine, N.J., Wilf, H.S.: Uniqueness theorems for periodic functions. *Proceedings of the American*  
495 *Mathematical Society* **16**(1), 109–114 (1965)
- 496 14 Gagie, T., Gawrychowski, P., Kärkkäinen, J., Nekrich, Y., Puglisi, S.J.: A faster grammar-  
497 based self-index. In: *Proc. 6th International Conference on Language and Automata Theory*  
498 *and Applications (LATA)*. pp. 240–251. LNCS 7183 (2012)
- 499 15 Gagie, T., Gawrychowski, P., Kärkkäinen, J., Nekrich, Y., Puglisi, S.J.: LZ77-based self-  
500 indexing with faster pattern matching. In: *Proc. 11th Latin American Symposium on Theoret-*  
501 *ical Informatics (LATIN)*. pp. 731–742 (2014)
- 502 16 Gagie, T., Navarro, G., Prezza, N.: Fully-functional suffix trees and optimal text searching in  
503 BWT-runs bounded space. *Journal of the ACM* **67**(1), article 2 (2020)
- 504 17 Ganardi, M., Jez, A., Lohrey, M.: Balancing straight-line programs. *Journal of the ACM*  
505 **68**(4), 27:1–27:40 (2021)
- 506 18 Gao, Y.: Computing matching statistics on repetitive texts. In: *Proc. 32nd Data Compression*  
507 *Conference (DCC)*. pp. 73–82 (2022)
- 508 19 Gawrychowski, P., Karczmarz, A., Kociumaka, T., Lacki, J., Sankowski, P.: Optimal dynamic  
509 strings. In: *Proc. 29th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*. pp.  
510 1509–1528 (2018)
- 511 20 Jez, A.: Approximation of grammar-based compression via recompression. *Theoretical Com-*  
512 *puter Science* **592**, 115–134 (2015)
- 513 21 Jez, A.: A really simple approximation of smallest grammar. *Theoretical Computer Science*  
514 **616**, 141–150 (2016)
- 515 22 Kärkkäinen, J., Ukkonen, E.: Lempel-Ziv parsing and sublinear-size index structures for string  
516 matching. In: *Proc. 3rd South American Workshop on String Processing (WSP)*. pp. 141–155  
517 (1996)
- 518 23 Karp, R.M., Rabin, M.O.: Efficient randomized pattern-matching algorithms. *IBM Journal of*  
519 *Research and Development* **2**, 249–260 (1987)

- 520 24 Kempa, D., Prezza, N.: At the roots of dictionary compression: String attractors. In: Proc.  
521 50th Annual ACM Symposium on the Theory of Computing (STOC). pp. 827–840 (2018)
- 522 25 Kempa, D., Kociumaka, T.: Collapsing the hierarchy of compressed data structures: Suffix  
523 arrays in optimal compressed space. In: Proc. 64th IEEE Annual Symposium on Foundations  
524 of Computer Science (FOCS). pp. 1877–1886 (2023)
- 525 26 Kieffer, J.C., Yang, E.H.: Grammar-based codes: A new class of universal lossless source  
526 codes. *IEEE Transactions on Information Theory* **46**(3), 737–754 (2000)
- 527 27 Kociumaka, T., Navarro, G., Olivares, F.: Near-optimal search time in  $\delta$ -optimal space, and  
528 vice versa. *Algorithmica* **86**(4), 1031–1056 (2024)
- 529 28 Kociumaka, T., Navarro, G., Prezza, N.: Toward a definitive compressibility measure for  
530 repetitive sequences. *IEEE Transactions on Information Theory* **69**(4), 2074–2092 (2023)
- 531 29 Kociumaka, T., Radoszewski, J., Rytter, W., Walen, T.: Internal pattern matching queries in a  
532 text and applications. In: Proc. 26th Annual ACM-SIAM Symposium on Discrete Algorithms  
533 (SODA). pp. 532–551 (2015)
- 534 30 Kreft, S., Navarro, G.: On compressing and indexing repetitive sequences. *Theoretical Com-  
535 puter Science* **483**, 115–133 (2013)
- 536 31 Larsson, J., Moffat, A.: Off-line dictionary-based compression. *Proceedings of the IEEE* **88**(11),  
537 1722–1732 (2000)
- 538 32 Lempel, A., Ziv, J.: On the complexity of finite sequences. *IEEE Transactions on Information  
539 Theory* **22**(1), 75–81 (1976)
- 540 33 Maruyama, S., Sakamoto, H., Takeda, M.: An online algorithm for lightweight grammar-based  
541 compression. *Algorithms* **5**(2), 214–235 (2012)
- 542 34 Navarro, G.: Spaces, trees and colors: The algorithmic landscape of document retrieval on  
543 sequences. *ACM Computing Surveys* **46**(4), article 52 (2014), 47 pages
- 544 35 Navarro, G.: Indexing highly repetitive string collections, part I: Repetitiveness measures.  
545 *ACM Computing Surveys* **54**(2), article 29 (2021)
- 546 36 Navarro, G.: Indexing highly repetitive string collections, part II: Compressed indexes. *ACM  
547 Computing Surveys* **54**(2), article 26 (2021)
- 548 37 Navarro, G.: Computing MEMs on repetitive text collections. In: Proc. 34th Annual Sym-  
549 posium on Combinatorial Pattern Matching (CPM). p. article 22 (2023)
- 550 38 Navarro, G., Olivares, F., Urbina, C.: Balancing run-length straight-line programs. In: Proc.  
551 29th International Symposium on String Processing and Information Retrieval (SPIRE). pp.  
552 117–131 (2022)
- 553 39 Navarro, G., Prezza, N.: Universal compressed text indexing. *Theoretical Computer Science*  
554 **762**, 41–50 (2019)
- 555 40 Navarro, G.: Document listing on repetitive collections with guaranteed performance. *Theor-  
556 etical Computer Science* **772**, 58–72 (2019)
- 557 41 Navarro, G.: Computing MEMs and relatives on repetitive text collections. *ACM Transactions  
558 on Algorithms* **21**(1), article 12 (2025)
- 559 42 Nevill-Manning, C., Witten, I., Mautsby, D.: Compression by induction of hierarchical  
560 grammars. In: Proc. 4th Data Compression Conference (DCC). pp. 244–253 (1994)
- 561 43 Nishimoto, T., I, T., Inenaga, S., Bannai, H., Takeda, M.: Fully dynamic data structure for  
562 LCE queries in compressed space. In: Proc. 41st International Symposium on Mathematical  
563 Foundations of Computer Science (MFCS). pp. 72:1–72:15 (2016)
- 564 44 Raskhodnikova, S., Ron, D., Rubinfeld, R., Smith, A.: Sublinear algorithms for approximating  
565 string compressibility. *Algorithmica* **65**, 685–709 (2013)
- 566 45 Rytter, W.: Application of Lempel-Ziv factorization to the approximation of grammar-based  
567 compression. *Theoretical Computer Science* **302**(1-3), 211–222 (2003)
- 568 46 Sakamoto, H.: A fully linear-time approximation algorithm for grammar-based compression.  
569 *Journal of Discrete Algorithms* **3**(2–4), 416–430 (2005)
- 570 47 Storer, J.A., Szymanski, T.G.: Data compression via textual substitution. *Journal of the ACM*  
571 **29**(4), 928–951 (1982)

## **XX:16** Counting on General Run-Length Grammars

- 572 **48** Tsuruta, K., Köppl, D., Nakashima, Y., Inenaga, S., Bannai, H., Takeda, M.: Grammar-  
573 compressed self-index with Lyndon words. CoRR **2004.05309** (2020)