

# Language Operations with Regular Expressions of Polynomial Size<sup>1</sup>

Hermann Gruber<sup>a,\*</sup>,<sup>2</sup> Markus Holzer<sup>a,2</sup>

<sup>a</sup>*Institut für Informatik, Justus-Liebig-Universität Giessen, Arndtstraße 2,  
D-35392 Giessen, Germany*

---

## Abstract

This work deals with questions regarding to what extent regularity-preserving language operations affect the descriptive complexity of regular expressions. Some language operations are identified which are feasible for regular expressions in the sense that the result of the operation can be represented as a regular expression of size polynomial in that of the operands. We prove that taking language quotients, in particular the prefix and suffix closures, of a regular set can incur at most a quadratic blow-up on the required expression size. The circular shift operation can cause only a cubic increase in size and an at least quadratic bloat can be necessary in the worst case.

*Key words:* regular expressions, derivatives, language quotient, cyclic shift, circular shift

---

## 1 Introduction

In the last 20 years, a large body of research on the descriptive complexity of finite automata has been developed. To the authors' knowledge, the first

---

\* Corresponding author.

*Email addresses:* [hermann.k.gruber@informatik.uni-giessen.de](mailto:hermann.k.gruber@informatik.uni-giessen.de) (Hermann Gruber), [holzer@informatik.uni-giessen.de](mailto:holzer@informatik.uni-giessen.de) (Markus Holzer).

<sup>1</sup> This paper is a completely revised and expanded version of a paper presented at the 10th Workshop on Descriptive Complexity of Formal Systems (DCFS) held in Charlottetown, Prince Edward Island, Canada, July 16–18, 2008.

<sup>2</sup> Most of the work was done while the first author was at Institut für Informatik, Ludwig-Maximilians-Universität München, Oettingenstraße 67, D-80538 München, Germany, and the second author was at Institut für Informatik, Technische Universität München, Boltzmannstraße 3, D-85748 Garching bei München, Germany.

systematic attempt to start a parallel development for the descriptive complexity of regular expressions was presented by Ellul et al. [9] at the workshop on “Descriptive Complexity of Formal Systems” (DCFS), in 2002. In particular, they raised the question of determining how basic language operations such as complementation and intersection affect the required regular expression size. Yet it is worth mentioning that a rather special language operation arising in the context of SGML specifications, namely that of removing the empty word, had been studied already some years before: It was shown in [23] that this operation can incur at most a quasilinear increase in regular expression size. In contrast, for the intersection and shuffle operation, exponential lower bounds have been shown recently, and complementation can even incur a doubly-exponential blow-up [11,13]. In [13] it was shown that the star height of a regular language is at most logarithmic in the minimum regular expression size, and lower bounds are proved by finding families of languages for which the respective language operations give rise to a dramatic increase in star height. In contrast, it is well known that taking language quotients does not increase the star height [7]. This and similar language operations appear to be a natural testing ground for deepening our understanding of the descriptive complexity of regular expressions: Either one has to find some new lower bound techniques, or one has to find a nontrivial implementation of these operations on regular expressions, or both—a straightforward procedure would be to convert the expression into a finite automaton, implement the operation on a finite automaton, and convert back to a regular expression using state elimination. Yet that last step can incur an exponential blow-up in general, even over binary alphabets [13].

Here, we give polynomial upper bounds for the required expression size resulting from taking language quotients and circular shift. Descriptive complexity aspects of these operations were already studied in [1,22] for the circular shift and [16,20] for language quotients—the latter two references consider deterministic finite automata with multiple start states, but the results easily translate to state complexity results for (left) quotients. The basic idea is to implement the operation for the special case of linear expressions [2], called single-occurrence regular expressions in [11]. These are expressions in which every alphabetic symbol occurs exactly once, which makes it easier to deal with as they can describe only local languages. To cover the general case, we study the interplay of the operations with length-preserving homomorphisms.

## 2 Basic Definitions

We recall some basic notions in formal language theory—for a thorough treatment, the reader might want to consult a textbook such as [18]. In particular, let  $\Sigma$  be a finite alphabet and  $\Sigma^*$  the set of all words over the alphabet  $\Sigma$ , in-

cluding the empty word  $\lambda$ . A (*formal*) *language* over the alphabet  $\Sigma$  is a subset of  $\Sigma^*$ . Apart from the regular operations on languages, namely (finite) union, catenation, and star, we briefly recall the following operations on languages: The *reversal* of a language  $L$ , denoted by  $L^R$ , consists of all words which, when read backwards yield a word in  $L$ . The (*left*) *derivative* of a language  $L$  with respect to a word  $w$ , written as  $w^{-1}L$ , is defined as  $\{x \mid wx \in L\}$ , the (*left*) *quotient* of  $L$  with respect to a set of words  $W$ , denoted by  $W^{-1}L$ , is defined as  $\bigcup_{w \in W} w^{-1}L$ . The special case  $W = \Sigma^*$  is known as the *suffix closure* of  $L$  and denoted by  $\text{suf}(L)$ . We can perform similar operations when reading words from right to left: The *right derivative* of a language  $L$  with respect to a word  $w$  is defined as  $\{v \mid vw \in L\}$ . This operation can be expressed using derivatives and reversal as  $((w^R)^{-1}L^R)^R$ ; right quotients and the prefix closure  $\text{pre}(L)$  are defined in an analogous manner. The circular (or cyclic) shift of a language, denoted by  $\circlearrowleft(L)$ , is given by  $\{xw \mid wx \in L\}$ .

Let  $\Sigma$  be an alphabet. The regular expressions over  $\Sigma$  are defined recursively in the usual way:<sup>3</sup>  $\emptyset$ ,  $\lambda$ , and every letter  $a$  with  $a \in \Sigma$  is a regular expression; and when  $s$  and  $t$  are regular expressions, then  $(s+t)$ ,  $(s \cdot t)$ , and  $(s)^*$  are also regular expressions. The language denoted by a regular expression  $r$ , denoted by  $L(r)$ , is defined as follows:  $L(\emptyset) = \emptyset$ ,  $L(\lambda) = \{\lambda\}$ ,  $L(a) = \{a\}$ ,  $L(s+t) = L(s) \cup L(t)$ ,  $L(s \cdot t) = L(s) \cdot L(t)$ , and  $L(s^*) = L(s)^*$ . Two regular expressions are called *equivalent* if they denote the same language. For a regular expression  $r$ , define  $\lambda(r) = \lambda$  if  $\lambda \in L(r)$ , and  $\lambda(r) = \emptyset$  otherwise. Likewise, for a language  $L$ , we define  $\lambda(L)$  analogously.

The *size* or *alphabetic width* of a regular expression  $r$  over the alphabet  $\Sigma$ , denoted by  $\text{alph}(r)$ , is defined as the total number of occurrences of letters of  $\Sigma$  in  $r$ . For a regular language  $L$ , we define its alphabetic width,  $\text{alph}(L)$ , as the minimum alphabetic width among all regular expressions describing  $L$ . The *star height* of a regular expression  $r$ , denoted by  $h(r)$  is a structural complexity measure inductively defined by

- (1)  $h(r) = 0$ , for  $r \in \Sigma \cup \{\emptyset, \lambda\}$ ,
- (2)  $h(s \cdot t) = h(s + t) = \max(h(s), h(t))$ , and
- (3)  $h(r^*) = 1 + h(r)$ .

The star height of a regular language  $L$  is then defined as the minimum star height among all regular expressions describing  $L$ .

For illustration of these concepts, the alphabetic width of the expression  $\lambda + ((ab)^*b)^*a^*b$  is 5, and its star height is 2. These numbers are also immediate

---

<sup>3</sup> For convenience, parentheses in regular expressions are sometimes omitted and the concatenation is simply written as juxtaposition. The priority of operators is specified in the usual fashion: Concatenation is performed before union, and star before both product and union.

upper bounds on the alphabetic width and star height, respectively, of the described language.

Let  $r$  be a regular expression. Following [19], we say that  $r$  is *reduced* if all of the following conditions hold: If  $r$  contains the symbol  $\emptyset$ , then  $r = \emptyset$ ; the expression  $r$  contains no subexpression of the form  $st$  or  $ts$ , satisfying  $L(s) = \{\lambda\}$  and no subexpression of the form  $(s^*)^*$ ; if  $r$  contains a subexpression of the form  $s+t$  or  $t+s$  with  $L(s) = \{\lambda\}$ , then  $\lambda \notin L(t)$ ; if  $r$  contains a subexpression of the form  $s^*$ , then  $L(s) \neq \{\lambda\}$ . Otherwise  $r$  is called reducible. The above definition suggests some rewriting rules, such as replacing  $s + \emptyset$  with  $s$ , and a few more rules, see [19]. By iteratively applying the rules to all subexpressions until none is applicable, we can *reduce* every regular expression to arrive at a reduced one.

Clearly, for every regular expression there exists an equivalent reduced regular expression with alphabetic width and star height no larger than the original expression. We will need the following relation between star height and alphabetic width of reduced regular expressions:

**Lemma 1** *Let  $r$  be a reduced regular expression. Then  $h(r) \leq \text{alph}(r)$ .*

**PROOF.** First, consider the cases  $L(r) = \emptyset$  and  $L(r) = \{\lambda\}$ . It is easy to see from the definition of reduced expressions that we must have  $r = \emptyset$  and  $r = \lambda$ , respectively, and the claim holds in these cases. These are the only reduced expressions with alphabetic width 0.

For the case  $\text{alph}(r) \geq 1$ , we prove the following claim by induction on the total number of occurrences of operators in  $r$ : If the reduced expression  $r$  is a starred expression, then  $h(r) \leq \text{alph}(r)$ , otherwise  $h(r) \leq \text{alph}(r) - 1$ .

If  $r$  contains no operators at all, then the statement clearly holds. To do the induction step, assume the statement holds for all regular expressions with at most  $m$  occurrences of operators. If  $r$  is of the form  $(s)^*$ , then  $s$  is not a starred expression, since  $r$  is reduced. Furthermore,  $L(s)$  can be neither empty nor equal to  $\{\lambda\}$ , since otherwise  $L(r) = \{\lambda\}$ , and  $r$  would not be reduced. Thus by induction assumption,  $h(s) \leq \text{alph}(s) - 1$ . Since  $\text{alph}(r) = \text{alph}(s)$  and  $h(r) = h(s) + 1$ , the claim also holds in this case. In the cases  $r = s + t$  and  $r = s \cdot t$ , we have  $h(r) = \max(h(s), h(t))$ . The total number of occurrences of operators in  $r$  and  $s$  is each strictly less than the number of occurrences in  $r$ , so the induction hypothesis applies. However, the induction does not go through directly in the subcase where both  $s$  and  $t$  are stars, that is,  $s = \hat{s}^*$  and  $t = \hat{t}^*$ . Nevertheless, since both  $s$  and  $t$  are reduced, neither  $\hat{s}$  nor  $\hat{t}$  are stars, and we can apply the induction hypothesis to them and obtain

$$h(\hat{s}) + h(\hat{t}) \leq \text{alph}(\hat{s}) + \text{alph}(\hat{t}) - 2.$$

Regarding the left-hand side of this inequality, note that

$$h(r) - 1 = \max\{h(\hat{s}), h(\hat{t})\} \leq h(\hat{s}) + h(\hat{t}).$$

Regarding the right-hand side, note that  $\text{alph}(\hat{s}) + \text{alph}(\hat{t}) = \text{alph}(r)$ . This completes the induction step also in this subcase, and the lemma is proved.  $\square$

### 3 Linear Expressions

Let  $r$  be a regular expression over the alphabet  $\Sigma$ . Recall that the alphabetic width of  $r$ , denoted by  $\text{alph}(r)$ , is the total number of occurrences of symbols in  $r$ . We refer to the  $i$ th alphabetic letter in  $r$  as the  $i$ th *position*. A regular expression  $r$  over an alphabet  $\Sigma = \{a_1, a_2, \dots, a_n\}$  is called a *linear expression* if and only if  $|\Sigma| = \text{alph}(r)$  and the  $i$ th position in  $r$  is the symbol  $a_i$ . In this case, there is a straightforward bijection between positions and alphabet symbols, and here we shall often denote the used alphabet by  $P_r$ .

For two alphabets  $\Sigma$  and  $\Gamma$ , a homomorphism  $h : \Sigma^* \rightarrow \Gamma^*$  is *length-preserving* or also *letter-to-letter* if it maps all symbols from  $\Gamma$  to symbols from  $\Sigma$ . It is easy to see that each regular expression  $r$  is the image of a unique linear expression  $\bar{r}$  under a length-preserving homomorphism: That homomorphism maps the symbol  $a_i$  to the  $i$ th position of  $r$ . This homomorphism will be denoted by  $\ell_r$  or just  $\ell$  in the case  $r$  is understood from the context.

**Example 2** For the regular expression  $r = ((ab)^*a)^*$ , the corresponding linear expression is  $\bar{r} = ((a_1a_2)^*a_3)^*$ , and the length-preserving homomorphism which maps  $\bar{r}$  to  $r$  is given by  $\ell_r = \{a_1 \mapsto a, a_2 \mapsto b, a_3 \mapsto a\}$ .  $\square$

Let  $\Sigma$  be an alphabet. A language  $L \subseteq \Sigma^*$  is *local* if it can be written as

$$L = \lambda(L) \cup (P\Sigma^* \cap \Sigma^*S) \setminus (\Sigma^*N\Sigma^*),$$

for some  $P, S \subseteq \Sigma$  and  $N \subseteq \Sigma^2$ . Note that in this definition, we permit the empty word to be a member of a local language. The concept of local languages is related to linear expressions as follows [3]:

**Theorem 3** For every linear expression  $r$ , the language  $L(r)$  is local.

We briefly recall the definition of the canonical derivative  $d_a(r)$  of a linear expression  $r$  with respect to an alphabet symbol  $a$ , in the reformulation given in [6, Prop. 6]:

**Definition 4** Let  $r$  be a reduced linear expression and let  $a$  be a symbol in  $P_r$ .

Then the canonical derivative  $d_a(r)$  is computed recursively by applying the following rules and finally reducing the expression:

$$\begin{aligned}
d_a(a) &= \lambda \\
d_a(s + t) &= \begin{cases} d_a(s) & \text{if } d_a(s) \neq \emptyset \\ d_a(t) & \text{otherwise} \end{cases} \\
d_a(s \cdot t) &= \begin{cases} d_a(s) \cdot t & \text{if } d_a(s) \neq \emptyset \\ d_a(t) & \text{otherwise} \end{cases} \\
d_a(s^*) &= d_a(s) \cdot s^*
\end{aligned}$$

And  $d_a(r) = \emptyset$  in all cases not covered above.

The study [6] relates the canonical derivatives of a linear expression to the original definition of derivatives for general regular expressions due to Brzozowski [5], and to the continuations introduced by Berry and Sethi [2]. The results from [6] relevant to our context are summarized in the following characterization:

**Theorem 5** *Let  $r$  be a linear expression, let  $a$  be a symbol in  $P_r$ , and  $u$  a word over  $P_r$ . If the set  $(ua)^{-1}L(r)$  is nonempty, then it is described by the canonical derivative  $d_a(r)$ .*

Thus for a reduced linear expression  $r$ , the canonical derivative  $d_a(r)$  describes the language quotient  $(P_r^*a)^{-1}L(r)$ .

**Example 6** *Consider again the linear expression  $\bar{r} = ((ab)^*c)^*$  from Example 2; we now use the alphabet  $\{a, b, c\}$  instead of  $\{a_1, a_2, a_3\}$  to increase readability. Then*

$$\begin{aligned}
d_a(\bar{r}) &= d_a((ab)^*c) \cdot ((ab)^*c)^* = d_a((ab)^*) \cdot c \cdot ((ab)^*c)^* \\
&= d_a(ab) \cdot (ab)^* \cdot c \cdot ((ab)^*c)^* = d_a(a) \cdot b \cdot (ab)^* \cdot c \cdot ((ab)^*c)^* \\
&= \lambda \cdot b \cdot (ab)^* \cdot c \cdot ((ab)^*c)^* = b(ab)^*c((ab)^*c)^*.
\end{aligned}$$

A similar computation yields  $d_b(\bar{r}) = (ab)^*c((ab)^*c)^*$  and  $d_c(\bar{r}) = ((ab)^*c)^*$ . □

Now we generalize the above notion from symbols  $a \in P_r$  to sets of symbols  $A \subseteq P_r$  as follows:

**Definition 7** *Let  $r$  be a reduced linear expression and let  $A$  be a set of symbol in  $P_r$ . Then the canonical derivative  $d_A(r)$  is computed recursively by applying*

the following rules and finally reducing the expression:

$$\begin{aligned}
d_A(a) &= \lambda \quad \text{if } a \in A \\
d_A(s + t) &= d_B(s) + d_{A \setminus B}(t) \\
d_A(s \cdot t) &= d_B(s) \cdot t + d_{A \setminus B}(t) \\
d_A(s^*) &= d_A(s) \cdot s^*
\end{aligned}$$

with  $B = \{a \in A \mid d_a(s) \neq \emptyset\}$ . And  $d_A(r) = \emptyset$  in all cases not covered above.

A straightforward structural induction shows that the definition works as expected:

**Lemma 8** *Let  $r$  be a linear expression and  $A$  a set of symbols in  $r$ . Then  $L(d_A(r)) = \bigcup_{a \in A} L(d_a(r))$ .  $\square$*

**Example 9** *Again we consider the linear expression  $\bar{r} = ((ab)^*c)^*$ . Let  $A = \{b, c\}$ , then the expression  $d_A(\bar{r})$  computes as follows:*

$$\begin{aligned}
d_A(\bar{r}) &= d_A((ab)^*c) \cdot ((ab)^*c)^* = (d_b((ab)^*) \cdot c + d_c(c)) \cdot ((ab)^*c)^* \\
&= (d_b(ab) \cdot (ab)^* \cdot c + d_c(c)) \cdot ((ab)^*c)^* \\
&= (d_b(b) \cdot (ab)^* \cdot c + d_c(c)) \cdot ((ab)^*c)^* = (\lambda \cdot (ab)^* \cdot c + \lambda) \cdot ((ab)^*c)^* \\
&= ((ab)^*c + \lambda) \cdot ((ab)^*c)^*.
\end{aligned}$$

In the last line of the above computation, we applied two different rules for reducing the expression. Note that this expression is indeed shorter than the expression  $d_b(\bar{r}) + d_c(\bar{r})$  we computed in Example 6, although one might first be tempted to expect the contrary.

A similar computation yields the equalities  $d_{\{a,b\}}(\bar{r}) = ((b+\lambda)(ab)^*c) \cdot ((ab)^*c)^*$ ,  $d_{\{a,c\}}(\bar{r}) = (b(ab)^*c + \lambda) \cdot ((ab)^*c)^*$ , and last but not least we find  $d_{\{a,b,c\}}(\bar{r}) = ((b+\lambda)(ab)^*c + \lambda) \cdot ((ab)^*c)^*$ .  $\square$

Next, we estimate the size of the expressions  $d_A(r)$ .

**Lemma 10** *Let  $r$  be a reduced linear expression of alphabetic width  $n \geq 1$  and star height  $h$ , and let  $A$  be a subset of  $P_r$ . Then the expression  $d_A(r)$  has size at most  $\frac{n^2-n}{2} + hn = O(n^2)$ .*

**PROOF.** First, recall that Lemma 1 implies that the claimed size is in  $O(n^2)$ . We prove the claim by induction on the depth  $d \geq 0$  of the syntax tree of  $r$ . In the case  $d = 0$ , then with  $\text{alph}(r) \geq 1$  we must have  $r = a$  for some  $a \in P_r$ , and the claim clearly holds. To do the induction step, we consider three cases.

If  $r$  is of the form  $s + t$ , then  $d_A(r)$  is the expression obtained from reducing  $d_B(s) + d_{A \setminus B}(t)$ . Let  $\text{alph}(s) = k$  and  $\text{alph}(t) = n - k$ , for some  $k \geq 0$ . By

induction hypothesis we obtain that

$$\text{alph}(d_A(r)) \leq \frac{k^2 - k}{2} + hk + \frac{(n - k)^2 - (n - k)}{2} + h(n - k).$$

By rearranging terms, we get

$$\frac{k^2 - k + (n - k)^2 - (n - k)}{2} = \frac{n^2 - n}{2} + k(k - n) \leq \frac{n^2 - n}{2},$$

and thus  $\text{alph}(d_A(r))$  is bounded above by  $\frac{n^2 - n}{2} + hn$  in this case.

If  $r$  is of the form  $s \cdot t$ , then  $d_A(r)$  is the expression obtained from reducing the expression  $d_B(s) \cdot t + d_{A \setminus B}(t)$ . Since  $r$  is reduced, both  $s$  and  $t$  have alphabetic width at least 1. Letting  $k \geq 1$  denote the alphabetic width of  $s$  and  $n - k$  the alphabetic width of  $t$ , we obtain by induction hypothesis that

$$\text{alph}(d_A(r)) \leq \frac{k^2 - k}{2} + hk + \frac{(n - k)^2 - (n - k)}{2} + h(n - k) + n - k.$$

By a similar computation as for the previous case, the right hand side in the above inequality is still bounded above by  $\frac{n^2 - n}{2} + hn$ .

Finally, if  $r$  is of the form  $s^*$ , then the depth and the star height of  $s$  are both smaller than those of  $r$ , and by induction assumption  $\text{alph}(d_A(r)) \leq \text{alph}(d_A(s)) + n \leq \frac{n^2 - n}{2} + (h - 1)n + n$ . This covers all possible cases, and the proof is completed.  $\square$

We remark that our notion of  $d_A(r)$  differs from the one given in [6] in that our definition yields expressions of size  $O(n^2)$ , while defining  $d_A$  as  $\sum_{a \in A} d_a(r)$  would be much more redundant. It should be said that the actual size of these expressions is immaterial in the context of [6], but is important in the present paper.

Now we take a closer look at local languages. The following lemma is easy to see from the definition of local languages:

**Lemma 11** *If  $L \subseteq \Sigma^*$  is a local language and  $a$  is a symbol in  $\Sigma$ , then*

- (1)  $u_1 \cdot a \cdot v_1 \in L$  and  $u_2 \cdot a \cdot v_2 \in L$  implies  $u_1 \cdot a \cdot v_2 \in L$ , and
- (2)  $ua \in \text{pre}(L)$  implies  $(ua)^{-1}L = (\Sigma^*a)^{-1}L$ .  $\square$

For a local language  $L$  and an alphabet symbol  $a$ , we may thus define  $d_a(L)$  as the language quotient  $d_a(L) = (\Sigma^*a)^{-1}L$ . Likewise, for a set of symbols  $A$  define  $d_A(L) = (\Sigma^*A)^{-1}L$ . This operator overloading is perfectly consistent with the use of the notation  $d_a(r)$  to denote the canonical derivative of a linear



expression  $r$  with respect to an alphabet symbol  $a$ : By Theorem 5, we have  $L(d_a(r)) = d_a(L(r))$ .

The above characterization allows us to provide a neat formula for left quotients of local languages:

**Lemma 12** *Let  $L \subseteq \Sigma^*$  be a local language and let  $W \subseteq \Sigma^*$  an arbitrary language. Define  $A = \{a \in \Sigma \mid W \cap \text{pre}(L) \cap \Sigma^*a \neq \emptyset\}$ . Then*

$$W^{-1}L = \lambda(W) \cdot L \cup \bigcup_{a \in A} d_a(L) = \lambda(W) \cdot L \cup d_A(L).$$

**PROOF.** First of all, it follows from the definition of language quotients that  $W^{-1}L = (W \cap \text{pre}(L))^{-1}L$ . Second, if we decompose the set  $W$  as  $\lambda(W) \cup \bigcup_{a \in \Sigma} (W \cap \Sigma^*a)$ , we obtain

$$\begin{aligned} W^{-1}L &= \lambda(W)^{-1} \cdot L \cup \bigcup_{a \in \Sigma} (W \cap \Sigma^*a \cap \text{pre}(L))^{-1}L \\ &= \lambda(W) \cdot L \cup \bigcup_{a \in A} (W \cap \Sigma^*a \cap \text{pre}(L))^{-1}L, \end{aligned} \tag{1}$$

since  $\bigcup_{a \in \Sigma \setminus A} (W \cap \Sigma^*a \cap \text{pre}(L))^{-1}L = (\emptyset)^{-1}L = \emptyset$ . Finally, for  $a \in A$ , let  $ua$  be any word in  $W \cap \Sigma^*a \cap \text{pre}(L)$ . Then from Lemma 11 one can readily deduce that  $(ua)^{-1}L = (\Sigma^*a \cap \text{pre}(L))^{-1}L = (\Sigma^*a)^{-1}L = d_a(L)$ . Thus  $(W \cap \Sigma^*a \cap \text{pre}(L))^{-1}L = d_a(L)$ . By putting this into Equation (1), the result follows.  $\square$

**Example 13** *We compute the left quotient of the language  $L = L(\bar{r})$  denoted by the linear expression  $\bar{r} = ((ab)^*c)^*$  with respect to the set  $W$  denoted by the regular expression  $\lambda + ((a+c)b)^*(a+c)$ . To this end, we identify the set  $A = \{d \in \{a, b, c\} \mid W \cap \Sigma^*d \cap \text{pre}(L) \neq \emptyset\}$ . No word in  $W$  ends with  $b$ , but the words  $a$  and  $c$  are both in  $W$  and prefixes of words in  $L$ , thus  $A = \{a, c\}$  in our case. In general, the set  $A$  can be effectively computed provided  $W$  is represented in a machine model that effectively allows intersection with regular sets and has a decidable emptiness problem. This is the case, e.g., for the finite automaton and pushdown automaton models, see [18].*

Now by Lemma 12 holds  $W^{-1}L = \lambda(W) \cdot L \cup d_{\{a,c\}}(L)$ . Here we have  $\lambda(W) \cdot L = L$ , and a regular expression denoting  $d_{\{a,c\}}(L)$  is given in the previous Example 9. Thus, a regular expression denoting the quotient  $W^{-1}L$  is given by  $\bar{r} + d_{\{a,c\}}(\bar{r}) = ((ab)^*c)^* + (b(ab)^*c + \lambda) \cdot ((ab)^*c)^*$ .  $\square$

Also for the circular shift of local languages, we obtain a nice characterization:

**Lemma 14** *Let  $L \subseteq \Sigma^*$  be a local language. Then for the circular shift  $\circlearrowleft(L)$  we have*

$$\circlearrowleft(L) = \lambda(L) \cup \bigcup_{a \in \Sigma} a \cdot d_a(L) \cdot (d_a(L^R))^R.$$

**PROOF.** The circular shift of any language  $L$  can by definition be written as  $\circlearrowleft(L) = \lambda(L) \cup \bigcup_{a \in \Sigma} L_a$ , where

$$L_a = \{ awv \mid v, w \in \Sigma^*, vaw \in L \}.$$

In particular, as  $L$  is local, Lemma 11 tells us that there cannot be any dependencies between the subwords  $v$  and  $w$  in the definition of  $L_a$ . Thus  $L_a$  can be rewritten as

$$L_a = \bigcup_{\{v \in \Sigma^* \mid va \cdot d_a(L) \subseteq L\}} a \cdot d_a(L) \cdot v,$$

and since local languages are readily seen to be closed under reversal, redoing the same trick for the reversed language yields  $L_a = a \cdot d_a(L) \cdot d_a(L^R)^R$ , as desired.  $\square$

**Example 15** *We also compute the circular shift of the language  $L = L(\bar{r})$  in our running example. By Lemma 14, we can write  $\circlearrowleft(L)$  as  $\lambda(L) \cup \bigcup_{a \in \Sigma} a \cdot d_a(L) \cdot (d_a(L^R))^R$ . The sets  $d_a(L)$ ,  $d_b(L)$  and  $d_c(L)$  are denoted by the expressions  $d_a(\bar{r}) = b(ab)^*c((ab)^*c)^*$ ,  $d_b(\bar{r}) = (ab)^*c((ab)^*c)^*$ , and  $d_c(\bar{r}) = ((ab)^*c)^*$ , respectively, which were computed in Example 6. In a similar manner, we obtain for  $\bar{r}^R = (c(ba)^*)^*$  the canonical derivatives  $d_a(\bar{r}^R) = (ba)^*(c(ba)^*)^*$ ,  $d_b(\bar{r}^R) = a(ba)^*(c(ba)^*)^*$ , and  $d_c(\bar{r}^R) = (ba)^*(c(ba)^*)^*$ . We have  $\lambda(L) = \lambda$ , and straightforward rules for implementing the reversal of regular expressions yield  $((ba)^*(c(ba)^*)^*)^R = ((ab)^*c)^*(ab)^*$  and  $(a(ba)^*(c(ba)^*)^*)^R = ((ab)^*c)^*(ab)^*a$ . Thus a regular expression denoting  $\circlearrowleft(L(\bar{r}))$  is given by*

$$\begin{aligned} & \lambda + a \cdot b(ab)^*c((ab)^*c)^* \cdot ((ab)^*c)^*(ab)^* \\ & + b \cdot (ab)^*c((ab)^*c)^* \cdot ((ab)^*c)^*(ab)^*a + c \cdot ((ab)^*c)^* \cdot ((ab)^*c)^*(ab)^*. \end{aligned}$$

*Finally, we note that in an actual implementation the computational overhead for the two reversal operations carried out here could be saved by defining the concept of “canonical right derivatives” using rules analogous to those for computing canonical (left) derivatives in Definition 4.  $\square$*

These characterizations immediately lend themselves to an implementation of quotient and circular shift operations on linear expressions via canonical derivatives. Using Lemma 10, we can estimate the resulting expression size as follows:

**Theorem 16** *Let  $r$  be a linear expression of size  $n$ , and let  $L = L(r)$ . Then for any set of words  $W \subseteq P_r^*$ , there is a regular expression of size  $O(n^2)$  denoting  $W^{-1}L$ , and a regular expression of size  $O(n^3)$  denoting the circular shift  $\circlearrowleft(L)$ .  $\square$*

## 4 The General Case

The above results allow us to compute from a given linear expression relatively small regular expressions denoting a language quotient or the circular shift of the denoted language. In this section, we investigate the interaction of (length-preserving) homomorphisms with the language operations under consideration to transfer the obtained results to the general case. The easier case is a language operation that commutes with length-preserving homomorphisms. This is the case for the circular shift:

**Lemma 17** *Let  $\ell$  be a length-preserving homomorphism, and let  $L \subseteq \Sigma^*$  be a language. Then  $\circ(\ell(L)) = \ell(\circ(L))$ .*

**PROOF.** Since both the homomorphism and circular shift operation commute with taking finite and infinite unions, it suffices to show the claim for the case where  $L$  contains a single word  $w = a_1 a_2 \dots a_k$ . In the case  $k \leq 1$ , we have  $L = \circ(L)$ , and the claim is trivially true. So assume  $k \geq 2$ . Then

$$\begin{aligned} \ell(\circ(\{w\})) &= \{\ell(w)\} \cup \ell(\{a_j \dots a_k a_1 a_2 \dots a_{j-1} \mid 2 \leq j \leq k\}) \\ &= \{\ell(w)\} \cup \{\ell(a_j) \dots \ell(a_k) \ell(a_1) \ell(a_2) \dots \ell(a_{j-1}) \mid 2 \leq j \leq k\}. \end{aligned}$$

Now let  $x = b_1 b_2 \dots b_k$ , with  $b_i = \ell(a_i)$ . Then

$$\begin{aligned} \circ(\{\ell(w)\}) &= \{\ell(w)\} \cup \{b_j \dots b_k b_1 b_2 \dots b_{j-1} \mid 2 \leq j \leq k\} \\ &= \{\ell(w)\} \cup \{\ell(a_j) \dots \ell(a_k) \ell(a_1) \ell(a_2) \dots \ell(a_{j-1}) \mid 2 \leq j \leq k\}, \end{aligned}$$

thus proving the desired equality.  $\square$

The next lemma shows how homomorphisms interact with taking left derivatives. We would like to warn the reader that in the following we use two very similar standard notations that should not be confused: The notion  $w^{-1}L$  denotes a derivative of the language  $L$ , whereas  $\ell^{-1}(w)$  is used to denote the preimage of  $w$  under the map  $\ell$ , that is,  $\ell^{-1}(w) = \{x \mid \ell(x) = w\}$ .

**Lemma 18** *Let  $L \subseteq \Sigma^*$  be a regular language, let  $\ell : \Sigma^* \rightarrow \Gamma^*$  be a length-preserving homomorphism, and let  $w \in \Gamma^*$ . Then  $w^{-1}\ell(L) = \bigcup_{x \in \ell^{-1}(w)} \ell(x^{-1}L)$ .*

**PROOF.** Let  $A = (Q, \Sigma, \delta, Q_0, F)$  be a nondeterministic finite automaton (possibly with multiple start states) accepting  $L$ , in the standard notation of [18]. We obtain a nondeterministic finite automaton  $B$  accepting  $\ell(L)$  by a standard construction: Let  $B = (Q, \Gamma, \delta', Q_0, F)$  with  $\delta'(q, a) = \bigcup_{b \in \ell^{-1}(a)} \delta(q, b)$ ,

for every  $q \in Q$  and every  $a \in \Gamma$ . For an automaton  $C$  accepting  $w^{-1}\ell(L)$ , we perform the standard quotient construction: Let  $C = (Q, \Gamma, \delta', Q'_0, F)$ , with  $Q'_0 = \bigcup_{q_0 \in Q_0} \delta'(q_0, w)$ .

For a finite automaton accepting the language  $\bigcup_{x \in \ell^{-1}(w)} x^{-1}L$ , we define  $D = (Q, \Sigma, \delta, Q'_0, F)$  with the same start states as above. Note that for the set  $Q'_0$  we have

$$Q'_0 = \bigcup_{q_0 \in Q_0} \delta'(q_0, w) = \bigcup_{x \in \ell^{-1}(w)} \bigcup_{q_0 \in Q_0} \delta(q_0, x),$$

so this automaton indeed accepts the quotient  $\bigcup_{x \in \ell^{-1}(w)} x^{-1}L$ . To get an automaton accepting the image under  $\ell$  of this language, we replace in  $D$ , similar to above, the transition function  $\delta$  with  $\delta'$  and change the input alphabet to  $\Gamma$ . But then we end up with the automaton  $(Q, \Gamma, \delta', Q'_0, F)$ , which is identical to the automaton  $C$ . Thus  $w^{-1}\ell(L) = \bigcup_{x \in \ell^{-1}(w)} \ell(x^{-1}L)$  as desired.  $\square$

Now we are ready to state the main result of this paper:

**Theorem 19** *Let  $r$  be a regular expression of size  $n$  denoting the language  $L \subseteq \Sigma^*$ , and let  $W \subseteq \Sigma^*$ . Then there is a regular expression of size  $O(n^2)$  denoting  $W^{-1}L$  and a regular expression of size  $O(n^3)$  denoting  $\circlearrowleft(L)$ .*

**PROOF.** Let  $\bar{r}$  be the linear expression for  $r$ , and  $\ell = \ell_r$  be the homomorphism which maps  $\bar{r}$  to  $r$ . Since  $\ell$  is length preserving, every word  $w \in \Sigma^*$  is in  $\ell(P_{\bar{r}})^*$ , and thus, by Lemma 18, we have

$$w^{-1}\ell(L(\bar{r})) = \bigcup_{x \in \ell^{-1}(w)} \ell(x^{-1}L(\bar{r})).$$

This readily generalizes to sets of words, and we obtain

$$\begin{aligned} W^{-1}\ell(L(\bar{r})) &= \bigcup_{w \in W} \bigcup_{x \in \ell^{-1}(w)} \ell(x^{-1}L(\bar{r})) \\ &= \ell\left(\bigcup_{x \in \ell^{-1}(W)} x^{-1}L(\bar{r})\right) = \ell\left((\ell^{-1}(W))^{-1}L(\bar{r})\right). \end{aligned}$$

The last one of the above expressions is the image under  $\ell$  of a quotient of  $\bar{r}$ . By Theorem 16, the latter language can be described by a regular expression of size  $O(n^2)$ , and applying the map  $\ell$  does not increase the alphabetic width. This shows that  $\text{alph}(W^{-1}L) = O(n^2)$ . For the circular shift, recall from Theorem 16 that  $\circlearrowleft(L(\bar{r}))$  has alphabetic width in  $O(n^3)$ . By Lemma 17,  $\circlearrowleft(L(\ell(\bar{r}))) = \ell(\circlearrowleft(L(\bar{r})))$ , and, as noted before, applying a length-preserving homomorphism does not increase the alphabetic width.  $\square$

**Example 20** *We compute the left-quotient of  $L(r)$  with respect to the set  $W$  denoted by the expression  $\lambda + (ab)^*a$ . Following the proof of Theorem 19, the*

language  $W^{-1}L(r)$  is equal to  $\ell((\ell^{-1}(W))^{-1}L(\bar{r}))$ , where  $\ell = \{a \mapsto a, b \mapsto b, c \mapsto a\}$  is the length-preserving homomorphism which maps  $\bar{r}$  to  $r$ . An expression denoting  $\ell^{-1}(W) = \ell^{-1}(\lambda + (ab)^*a)$  is given by  $\lambda + ((a+c)b)^*(a+c)$ . Recall from Example 13 that the quotient  $(\ell^{-1}(W))^{-1}L(\bar{r})$  is denoted by the expression

$$\bar{r} + d_{\{a,c\}}(\bar{r}) = ((ab)^*c)^* + (b(ab)^*c + \lambda) \cdot ((ab)^*c)^*.$$

Finally, applying the length-preserving homomorphism  $\ell$  to this expression yields the regular expression  $((ab)^*a) + (b(ab)^*a + \lambda)((ab)^*a)^*$  for  $W^{-1}L(r)$ .  $\square$

**Example 21** From Example 15 we deduce that the circular shift of  $L(r)$  is simply

$$\begin{aligned} \circlearrowleft(L(r)) &= \lambda + a \cdot b(ab)^*a((ab)^*a)^* \cdot ((ab)^*a)^*(ab)^* \\ &\quad + b \cdot (ab)^*a((ab)^*a)^* \cdot ((ab)^*a)^*(ab)^*a + a \cdot ((ab)^*a)^* \cdot ((ab)^*a)^*(ab)^*, \end{aligned}$$

by applying the length-preserving homomorphism  $\ell = \{a \mapsto a, b \mapsto b, c \mapsto a\}$  to the regular expression  $\circlearrowleft(L(\bar{r}))$ .  $\square$

Currently, we do not know whether these upper bounds have the right order of magnitude. At least, we can show a quadratic lower bound on the increase of alphabetic width for the circular shift operation. To this end, we make use of certain circular arrangements, which are well known in combinatorics: For an integer  $k \geq 1$ , a *binary De Bruijn sequence* of order  $k$  is a circular sequence  $\gamma$  of length  $2^k$  over the binary alphabet  $\{a, b\}$  such that every word in  $\{a, b\}^k$  is contained, at some *unique* position, in  $\gamma$ . The following theorem, proved first<sup>4</sup> in [10], is a classical result in combinatorics:

**Theorem 22** *For every  $k \geq 1$ , a binary De Bruijn sequence of order  $k$  exists.*

By cutting open a circular sequence  $\gamma$  at some specified position and read it clockwise, we obtain a word  $w \in \{a, b\}^*$  in the usual sense, of length  $|\gamma|$ . We call such a word a *cut* of  $\gamma$ . With these notions, we are now sufficiently armed to prove a lower bound on the alphabetic width on circular shift.

**Theorem 23** *There exist infinitely many regular languages  $L_n$  over a binary alphabet such that  $L_n$  admits a regular expression of alphabetic width  $n$ , but every regular expression describing  $\circlearrowleft(L_n)$  has alphabetic width at least  $\Omega(n^2)$ .*

**PROOF.** We take  $n = 2^k$ , for some  $k \geq 1$ . Let  $\gamma$  be a binary De Bruijn sequence of order  $k$ , and let  $w$  be a cut of  $\gamma$ . Our witness language is simply  $L_n = \{w\}$ . Of course,  $L_n$  admits a regular expression of alphabetic width at

<sup>4</sup> As acknowledged by De Bruijn [8], he later rediscovered essentially the same proof.

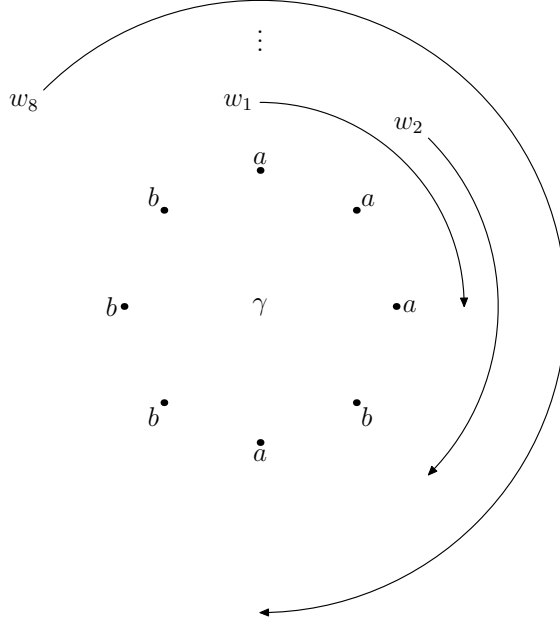


Fig. 1. The De Bruijn sequence  $\gamma = aaababbb$  of order  $k = 3$  and length  $n = 8$  is shown. The words  $w_i$ , for  $1 \leq i \leq n$ , used in the construction of the fooling set are obtained by cuts of  $\gamma$  at appropriate positions—here  $w_1 = aaababbb$  is obtained by cutting  $\gamma$  at 12-o'clock, while  $w_4 = babbbaaa$  is obtained by cutting at half past four.

most  $n$ . For a lower bound on the alphabetic width of  $\mathcal{C}(\{w\})$ , observe first that this set equals the set of all cuts of  $\gamma$ . Since there are  $n = 2^k$  such cuts, each of length  $n$ , there is trivially a regular expression of size  $n^2$  for this set. Our aim is to show that we cannot do essentially better.

Since every regular expression of size  $n$  can be transformed into a nondeterministic finite automaton having at most  $n+1$  states, it suffices to give a lower bound of  $\Omega(n^2)$  on the number of states required by any nondeterministic finite automaton accepting the set of all cuts of  $\gamma$ . To this end, we make use of the fact that the size of any fooling set for a regular language is a lower bound on the number of required states [4]. Here, a fooling set for a regular language  $L$  is a set of pairs of words  $(x_i, y_i)$  such that  $x_i y_i \in L$  for each  $i$ , but for  $i \neq j$ , *at least one* of the words  $x_i y_j$  and  $x_j y_i$  is not in  $L$ . For a detailed account on this and related lower bound techniques see, e.g., [12,17]. Now, let  $w_1 = w$  and for  $2 \leq i \leq n$ , we recursively define  $w_i = vu$ , if  $uv = w_{i-1}$ , where  $u \in \{a, b\}$  and  $v \in \Sigma^*$ —thus, all these words are obtained by cuts of  $\gamma$  at appropriate positions. For example, if we have  $w_1 = aaababbb$ , then  $w_4 = babbbaaa$ —see Figure 1.

Note that all these words are pairwise different because  $w$  is a cut of a De Bruijn sequence, i.e., for  $1 \leq i, j \leq n$  we have  $w_i \neq w_j$ , if  $i \neq j$ . For  $1 \leq i \leq n$  and  $1 \leq j \leq n$ , let  $x_{i,j}$  denote the prefix of length  $j$  of  $w_i$ . Similarly, we use  $y_{i,j}$

to denote the suffix of length  $j$  of  $w_i$ . We claim that

$$S = \{ (x_{i,j}, y_{i,n-j}) \mid 1 \leq i \leq n \text{ and } k \leq j \leq n - k \}$$

is a fooling set for the language  $\circ(\{w\})$  of all cuts of  $\gamma$ . Clearly, the word  $x_{i,j}y_{i,n-j}$  equals  $w_i$ , and every  $w_i$  is a cut of  $\gamma$ . We will prove that  $x_{i,j}y_{r,s}$  is not a member of  $\circ(\{w\})$ , if  $r \neq i$  or  $s \neq n - j$ . Consider first the case  $s \neq n - j$ . Since  $x_{i,j}$  is of length  $j$  and  $y_{r,s}$  of length  $s$ , their concatenation has length different from  $n$  and cannot be a cut. Now consider the case  $r \neq i$  and  $s = n - j$ . Recall that  $w_i$  is the *unique* cut of  $\gamma$  that begins with  $x_{i,k}$ . By rotational symmetry,  $w_i$  is also the *unique* cut that ends with  $y_{i,k}$ . For the sake of contradiction, assume now that  $x_{i,j}y_{r,n-j}$  is a cut of  $\gamma$ . Since the above word begins with  $x_{i,j}$  and  $j \geq k$ , it has  $x_{i,k}$  as prefix, so it must be equal to  $w_i$ . By similar means, since the word ends with  $y_{r,n-j}$  and  $n - j \geq k$ , it has  $y_{r,k}$  as suffix, so it must be equal to  $w_r$ . This implies  $w_i = w_r$ , a contradiction, since by assumption  $i \neq r$ .

Thus  $S$  is a fooling set for  $\circ(\{w\})$ , the set of all cuts of  $\gamma$ , and  $|S| = \Omega(n^2)$  is a lower bound for the number of states of any nondeterministic finite automaton accepting the language  $\circ(w)$ , as desired.  $\square$

Observe that this lower bound is tight if we restrict our attention to finite languages: If  $L(r)$  is finite, we can safely assume that  $r$  does not use the star operator, as observed in [15]: Any potential starred subexpression could describe only the empty word. In the case of expressions without star, it is readily seen from Definition 4 that all canonical derivatives are of linear size. Plugging this improved estimate into the proof of the  $O(n^3)$  bound for the circular shift in the general case gives a quadratic upper bound for the case of finite languages. We thus obtain a tight bound for the regular expression size of circular shift in the case of finite languages:

**Corollary 24** *Let  $r$  be a regular expression of size  $n$  denoting a finite language  $L \subset \Sigma^*$ . Then there is a regular expression of size  $O(n^2)$  denoting  $\circ(L)$ . In contrast, there exist infinitely many finite languages  $L_n$  over a binary alphabet such that  $L_n$  admits a regular expression of alphabetic width  $n$ , but every regular expression describing  $\circ(L_n)$  has alphabetic width at least  $\Omega(n^2)$ .*  $\square$

## 5 Conclusion

In this paper, we identified some regularity-preserving language operations whose effect on required regular expression size is not too drastic, i.e., which

can incur at most a polynomial blow-up. Among these are all operations which are special cases of language quotients, e.g., the prefix or suffix closure of a set of words, and the circular shift. The naive way to implement such an operation would involve a translation into finite automata and back. However, the conversion into the back direction likely causes an undesirable blow-up in expression size. In contrast, the algorithms presented here are entirely based on rewriting expressions and thus avoid these difficulties altogether. There are two ingredients in such an approach: First, the language operation under consideration needs to admit an efficient solution for linear expressions. Second, it needs to be somehow well-behaved with respect to length preserving homomorphisms.

One task for further research is to find other regularity preserving operations for which this or similar approaches might work. For instance, for the language of scattered substrings (superstrings, respectively) of the language described by a regular expression over  $\Sigma$ , we simply replace every position  $a$  with a subexpression  $\lambda + a$  (with a subexpression describing  $\Sigma^*a\Sigma^*$ , respectively) to obtain a regular expression denoting that language. Both operations can be thus performed with only linear increase in expression size provided  $\Sigma$  is fixed. Issues on the state complexity of these operations were studied recently in [14] and [21]. Another, probably difficult, challenge is to try to tighten the bounds given here. Quite a few lower bound techniques for regular expression size, apart from those based on the number of states required by a nondeterministic finite automaton, have been developed recently [11,13,15]. Apparently none of them can be used to infer something nontrivial about language quotients.

## References

- [1] P. R. J. Asveld. Generating all circular shifts by context-free grammars in Greibach normal form. *International Journal of Foundations of Computer Science*, 18(6):1139–1149, 2007.
- [2] G. Berry and R. Sethi. From regular expressions to deterministic automata. *Theoretical Computer Science*, 48(3):117–126, 1986.
- [3] J. Berstel and J. E. Pin. Local languages and the Berry-Sethi algorithm. *Theoretical Computer Science*, 155(2):439–446, 1996.
- [4] J.-C. Birget. Intersection and union of regular languages and state complexity. *Information Processing Letters*, 43:185–190, 1992.
- [5] J. A. Brzozowski. Derivatives of regular expressions. *Journal of the ACM*, 11(4):481–494, 1964.
- [6] J.-M. Champarnaud and D. Ziadi. Canonical derivatives, partial derivatives and finite automaton constructions. *Theoretical Computer Science*, 289(1):137–163, 2002.



- [7] R. S. Cohen and J. A. Brzozowski. General properties of star height of regular events. *Journal of Computer and System Sciences*, 4(3):260–280, 1970.
- [8] N. G. De Bruijn. Acknowledgement of priority to C. Flye Sainte-Marie on the counting of circular arrangements of  $2^n$  zeros and ones that show each  $n$ -letter word exactly once. T.H.-Report 75-WSK-06, Technological University Eindhoven, 1975.
- [9] K. Ellul, B. Krawetz, J. Shallit, and M. Wang. Regular expressions: New results and open problems. *Journal of Automata, Languages and Combinatorics*, 10(4):407–437, 2005.
- [10] C. Flye Sainte-Marie. Solution to question nr. 48. *L’Intermédiaire des mathématiciens*, 1:107–110, 1894. Reproduced as appendix in [8].
- [11] W. Gelade and F. Neven. Succinctness of the complement and intersection of regular expressions. In S. Albers and P. Weil, editors, *Proceedings of the 25th Symposium on Theoretical Aspects of Computer Science*, volume 08001 of *Dagstuhl Seminar Proceedings*, pages 325–336, 2008. Internationales Begegnungs- und Forschungszentrum fuer Informatik, Schloss Dagstuhl, Germany.
- [12] H. Gruber and M. Holzer. Finding lower bounds for nondeterministic state complexity is hard (extended abstract). In O. H. Ibarra and Z. Dang, editors, *Proceedings of the 10th International Conference on Developments in Language Theory*, number 4036 in LNCS, pages 363–374, Santa Barbara, California, USA, June 2006. Springer.
- [13] H. Gruber and M. Holzer. Finite automata, digraph connectivity, and regular expression size. In L. Aceto, I. Damgaard, L. A. Goldberg, M. M. Halldórsson, A. Ingólfssdóttir, and I. Walkuwiewicz, editors, *Proceedings of the 35th International Colloquium on Automata, Languages and Programming*, number 5126 of LNCS, pages 39–50, 2008. Springer.
- [14] H. Gruber, M. Holzer, and M. Kutrib. More on the size of Higman-Haines sets: Effective constructions. In J. O. Durand-Lose and M. Margenstern, editors, *Proceedings of the 5th International Conference Machines, Computations, and Universality*, number 4664 of LNCS, pages 193–204, 2008. Springer.
- [15] H. Gruber and J. Johannsen. Optimal lower bounds on regular expression size using communication complexity. In R. Amadio, editor, *Proceedings of the 11th International Conference on Foundations of Software Science and Computation Structures*, number 4962 of LNCS, pages 273–286, 2008. Springer.
- [16] M. Holzer, K. Salomaa, and S. Yu. On the state complexity of  $k$ -entry deterministic finite automata. *Journal of Automata, Languages and Combinatorics*, 6(4):453–466, 2001.
- [17] J. Hromkovič and G. Schnitger. On the hardness of determining small NFA’s and of proving lower bounds on their sizes. In M. Ito and F. M. Toyama, editors, *Proceedings of the 12th International Conference on Developments in Language Theory*, number 5257 of LNCS, pages 34–55, 2008. Springer.

- [18] J. E. Hopcroft and J. D. Ullman. *Introduction to Automata Theory, Languages and Computation*. Addison-Wesley, 1979.
- [19] L. Ilie and S. Yu. Follow automata. *Information and Computation*, 186(1):140–162, 2003.
- [20] M. Kappes. Descriptive complexity of deterministic finite automata with multiple initial states. *Journal of Automata, Languages and Combinatorics*, 5(3):269–278, 2000.
- [21] A. Okhotin. On the state complexity of scattered substrings and superstrings. TUCS Technical Report No.849, Turku Centre for Computer Science, 2007.
- [22] A. Okhotin and G. Jirásková. State complexity of cyclic shift. *RAIRO – Theoretical Informatics and Applications*, 42(2):335–360, 2008.
- [23] D. Ziadi. Regular expression for a language without empty word. *Theoretical Computer Science*, 163(1–2):309–315, 1996.