

# On the Representation of Block Languages

Guilherme Duarte<sup>1</sup>, Nelma Moreira<sup>1</sup>, Luca Prigioniero<sup>2</sup>, and Rogério Reis<sup>1\*</sup>

CMUP & DCC, Faculdade de Ciências da Universidade do Porto  
Rua do Campo Alegre, 4169-007 Porto, Portugal

{guilherme.duarte,nelma.moreira,rogerio.reis}@fc.up.pt

<sup>2</sup> Department of Computer Science, Loughborough University, UK  
L.Prigioniero@lboro.ac.uk

**Abstract.** In this paper we consider block languages, namely sets of words having the same length, and we propose a new representation for these languages. In particular, given an alphabet of size  $k$  and a length  $\ell$ , these languages can be represented by bitmaps of size  $k^\ell$ , in which each bit indicates whether the correspondent word, according to the lexicographical order, belongs to the language (bit equal to 1) or not (bit equal to 0). This representation turns out to be a good tool for the investigation of several properties of block languages, making proofs simpler and reasoning clearer. After showing how to convert bitmaps into minimal deterministic and nondeterministic finite automata, we use this representation as a tool to study the deterministic and nondeterministic state complexity of block languages, as well as the costs of basic operations on block languages, in terms of the sizes of the equivalent finite automata.

## 1 Introduction

In the area of formal languages and automata theory, the class of regular languages is one of the most investigated. Classical recognizers for this class are finite automata, in both deterministic and nondeterministic variants.

The capabilities of these machines of representing languages in a more or less succinct way have been widely studied in the area of *descriptive complexity*. In this context, the *size* of a model is measured in terms of number of symbols used to write down its description. In the specific case of finite automata, it is often considered the number of states as a measure of complexity.

In this area, the minimality of finite automata has been studied. For example, it is well known that, given a language  $L$ , the deterministic finite automaton of minimal size accepting it is unique (up to isomorphisms), and there exist efficient algorithms for the minimization of these machines [2]. The situation in the nondeterministic case is more challenging as a minimal nondeterministic finite automaton for  $L$  is not necessarily unique. Furthermore, given an integer  $n$ ,

---

\* This work was partially supported by CMUP, member of LASI, which is financed by national funds through FCT – Fundação para a Ciência e a Tecnologia, I.P., under the projects with reference UIDB/00144/2020 and UIDP/00144/2020.

deciding whether there is there a nondeterministic finite automaton with less than  $n$  states accepting  $L$  is a PSPACE-hard problem [15].

It is also classic to study the descriptonal complexity of operations; namely the cost, in terms of size, of the machines accepting the languages resulting from operations (e.g., union, intersection, complement, concatenation, reversal, etc.) on other devices [16,6].

In the literature, also several restrictions of well-known classes of languages are considered. The study of these specific cases are motivated by their applications and interesting because their costs are often lower than the general cases.

In this paper we consider finite languages where all words have the same length, which are called *homogeneous* or *block* languages. Their investigation is mainly motivated by their applications to several contexts such as code theory [10] and image processing [8,7]. As a subclass of finite languages, block languages inherit the size costs known for that class, which are known to be lower than the ones for the class of regular languages [4].

For instance, the minimisation of deterministic finite automata can be done in linear time in the case of finite (and hence also block) languages [14]. Due to the fact that all words have the same length, there are some gains in terms of descriptonal complexity. It is known that the elimination of nondeterminism from a  $n$ -state nondeterministic finite automaton for a block language costs  $2^{\Theta(\sqrt{n})}$  in size [8], which is smaller than the general case, for which the cost in size is  $2^{\Theta(n)}$  [12].

The maximum number of states of minimal deterministic finite automata for finite and block languages were studied by Câmpeanu and Ho [3], and Hanssen and Liu [9] determined the number of block languages that attain the maximum state complexity. Minimal deterministic finite automata for finite languages were enumerated by Almeida et al. [1]. Asymptotic estimates and exact formulae for the number of  $n$ -state minimal deterministic finite automata accepting finite languages over alphabets of size  $k$  were obtained by J. Priez [13] and by Price et. al [5]. On the other hand, the state complexity of operations over finite languages has been also well studied [6].

In this paper we propose a new representation for block languages. In particular, given an alphabet of size  $k$  and a length  $\ell$ , each block language can be represented by a binary string of length  $k^\ell$ , also called *bitmap*, in which each symbol (or *bit*) indicates whether the correspondent word, according to the lexicographical order, belongs to the language (bit equal to 1) or not (bit equal to 0). We use this representation as a tool to investigate several properties of block languages. More precisely, in Section 3 we show how to convert bitmaps into minimal deterministic and nondeterministic finite automata. While the former conversion can be done in polynomial time in the size of the bitmap, we prove the latter to be NP-complete. Then, in Section 4 we use this representation as a tool to study the deterministic state complexity of block languages, as well as the costs, in terms of the sizes of the equivalent finite automata, of basic operations on block languages.

## 2 Preliminaries

In this section we review some basic definitions about finite automata and languages and fix notation. Given an *alphabet* (finite set of *letters*)  $\Sigma$ , a *word*  $w$  is a sequence of symbols, and a *language*  $L \subset \Sigma^*$  is a set of words on  $\Sigma$ . The empty word is represented by  $\varepsilon$ . The (*left*) *quotient* of a language  $L$  by a word  $w \in \Sigma^*$  refers to the set  $w^{-1}L = \{w' \in \Sigma^* \mid ww' \in L\}$ . The *reversal* of  $L$ ,  $L^R$ , is the language obtained by reversing all the words in  $L$ .

A *nondeterministic finite automaton* (NFA) is a five-tuple  $\mathcal{A} = \langle Q, \Sigma, \delta, I, F \rangle$  where  $Q$  is a finite set of states,  $\Sigma$  is a finite alphabet,  $I \subseteq Q$  is the set of initial states,  $F \subseteq Q$  is the set of final states, and  $\delta : Q \times \Sigma \rightarrow 2^Q$  is the transition function. We consider the *size* of an NFA as its number of states. The transition function can be extended to words and to sets of states in the natural way. When  $I = \{q_0\}$ , we use  $I = q_0$ . We define the *finality* function  $\varepsilon$  on  $Q$  by  $\varepsilon(q) = \varepsilon$  if  $q \in F$  and  $\varepsilon(q) = \emptyset$ , otherwise. An NFA accepting a non-empty language is *trim* if every state is accessible from an initial state and every state leads to a final state. Given a state  $q \in Q$ , the *right language* of  $q$  is  $\mathcal{L}_q(\mathcal{A}) = \{w \in \Sigma^* \mid \delta(q, w) \cap F \neq \emptyset\}$ , and the *left language* is  $\overleftarrow{\mathcal{L}}_q(\mathcal{A}) = \{w \in \Sigma^* \mid q \in \delta(I, w)\}$ . The *language accepted* by  $\mathcal{A}$  is  $\mathcal{L}(\mathcal{A}) = \bigcup_{q \in I} \mathcal{L}_q(\mathcal{A})$ . An NFA  $\mathcal{A}$  is *minimal* if it has the smallest number of states among all NFAs that accept  $\mathcal{L}(\mathcal{A})$ . An NFA is *deterministic* (DFA) if  $|\delta(q, \sigma)| \leq 1$ , for all  $(q, \sigma) \in Q \times \Sigma$ , and  $|I| = 1$ . We can convert an NFA  $\mathcal{A}$  into an equivalent DFA  $D(\mathcal{A})$  using the well-known subset construction. Two states  $q_1, q_2$  are *equivalent* (or *indistinguishable*) if  $\varepsilon(q_1) = \varepsilon(q_2)$  and  $L_{q_1}(\mathcal{A}) = L_{q_2}(\mathcal{A})$ . If a DFA is *minimal* it has not equivalent states and it is unique up to renaming of states. If  $\mathcal{A}$  is a minimal DFA of  $L$ , then for each state  $q$ ,  $\mathcal{L}_q(\mathcal{A}) = w_q^{-1}L$  for some  $w_q \in \Sigma^*$  and if  $q \neq q'$  then  $w_q^{-1}L \neq w_{q'}^{-1}L$ . The *state complexity* of a language  $L$  is the size of its minimal DFA.

A trim NFA for a finite language  $\mathcal{A} = \langle Q, \Sigma, \delta, I, F \rangle$  is acyclic and ranked, i.e., the set of states  $Q$  is partitioned into  $Q_0 \cup Q_1 \cup \dots \cup Q_\ell$  such that every state  $q$  in rank  $Q_r$  reaches a final state by words of length  $r$  ( $Q_r = \{q \in Q \mid \exists w \in \Sigma^r \mid \delta(q, w) \in F\}$ ) and all transitions from states of rank  $Q_r$  lead only to states in  $Q_s$ , with  $s < r$  (no dead-state). We define the *width* of a rank  $r$  as the cardinality of the set  $Q_r$ , the *width* of  $\mathcal{A}$  to be the maximal width of a rank, i.e.,  $\text{width}(\mathcal{A}) = \max_{i \in [\ell]} |Q_i|$  and the *rank* of  $\mathcal{A}$ ,  $\text{rank}(\mathcal{A})$ , to be the length of its longest accepted word. A DFA for a finite language is also ranked but it may have a *dead-state*  $\Omega$  which is the only cyclic state, and is usually omitted. Formally  $\mathcal{A} = \langle Q \cup \{\Omega\}, \Sigma, \delta, q_0, F \rangle$  with  $F \subseteq Q$  and  $q_0 \neq \Omega$  such that  $(\forall \sigma \in \Sigma) \delta(\Omega, \sigma) = \Omega$  and  $(\forall x \in \Sigma^*)(\forall q \in Q) (\delta(q, x) \neq q)$ . In a trim acyclic automaton, two states  $q$  and  $q'$  are equivalent if they are both in the same rank, either final or not final, and the transition function is identical. An acyclic DFA can be minimised by merging equivalent states and the resulting algorithm runs in linear time in the size of the automaton (Revuz algorithm, [14,1]).

Given two integers  $i, j$  with  $i < j$ , let  $[i, j]$  denote the range from  $i$  to  $j$ , including both  $i$  and  $j$ , namely  $\{i, \dots, j\}$ . Moreover, we shall omit the left bound if it is equal to 0, thus  $[j] = \{0, \dots, j\}$ . In the following we shall consider sequences

of Boolean values,  $\mathbf{B} \in \{0, 1\}^n$  that we denote by *bitmaps*. Given two bitmaps  $\mathbf{B}_1, \mathbf{B}_2 \in \{0, 1\}^n$ ,  $\mathbf{B}_1 \circ \mathbf{B}_2$  represents the bitmap obtained by carrying out the bitwise operation  $\circ \in \{\vee, \wedge\}$ , and  $\overline{\mathbf{B}_1}$  the bitwise complement of  $\mathbf{B}_1$ .

### 3 Block Languages and Bitmaps

Given an alphabet  $\Sigma = \{\sigma_0, \dots, \sigma_{k-1}\}$  of size  $k > 0$  and an integer  $\ell \geq 0$ , a *block language*  $L \subseteq \Sigma^\ell$ , is a set of words of length  $\ell$  over  $\Sigma$ . The language  $L$  can be characterised by a word in  $\{0, 1\}^{k^\ell}$ , that we call *bitmap* denoted as

$$\mathbf{B}(L) = b_0 \cdots b_{k^\ell - 1},$$

where  $b_i = 1$  if the word  $w$  is in  $L$ , and  $i \in [k^\ell - 1]$  is the index of  $w$  in the lexicographical ordered list of all the words of  $\Sigma^\ell$ . In this case, we denote  $i$  by  $\text{ind}(w)$ . We will denote the bitmap of a language as  $\mathbf{B}$  when it is unambiguous to which language the bitmap refers to. Moreover, each bitmap  $\mathbf{B} \in \{0, 1\}^{k^\ell}$  represents a block language of length  $\ell$  over a  $k$ -ary alphabet, thus one can use any alphabet of size  $k$ .

*Example 1.* Let  $\Sigma = \{a, b\}$ ,  $\ell = 4$ , and

$$L = \{aaaa, aaba, aabb, abab, abba, abbb, babb, bbaa, bbab, bbba\}.$$

The bitmap of  $L$  is  $\mathbf{B}(L) = 1011011100011110$ .

A bitmap  $\mathbf{B} \in \{0, 1\}^{k^\ell}$  can be split into factors of length  $k^r$  for  $r \in [\ell]$ . Let  $s_i^r = b_{i \cdot k^r} \cdots b_{(i+1) \cdot k^r - 1}$  denote the  $i$ -th factor of length  $k^r$ , for  $i \in [k^{\ell-r} - 1]$ . Since each factor of length  $k^r$  can also be split into  $k$  factors,  $s_i^r$  is inductively defined as:

$$s_i^r = \begin{cases} b_i, & \text{if } r = 0, \\ s_{i \cdot k}^{r-1} \cdots s_{i \cdot k + k - 1}^{r-1}, & \text{otherwise.} \end{cases}$$

The following lemma formally introduces the observation that each factor  $s_i^r$ , with  $r \in [\ell]$  and  $i \in [k^{\ell-r} - 1]$ , represents a quotient of  $L$ .

**Lemma 1.** *Let  $L \subseteq \Sigma^\ell$  be a block language,  $|\Sigma| = k$ , and  $\mathbf{B}$  the bitmap of  $L$ . Let  $r \in [\ell]$ ,  $i \in [k^{\ell-r} - 1]$ , and  $w \in \Sigma^{\ell-r}$  be the word of index  $i$  of size  $\ell - r$ , in lexicographic order. Then, the factor  $s_i^r$  corresponds to the bitmap of  $w^{-1}L$ .*

*Example 2.* Recall the bitmap of Example 1,  $\mathbf{B} = 1011011100011110$ , with  $k = 2$  and  $\ell = 4$ . We have that  $s_0^2 = 1011$  refers to the bitmap of  $(aa)^{-1}L = \{aa, ba, bb\}$ ,  $s_1^3 = 00011110$  to the bitmap of  $b^{-1}L = \{abb, baa, bab, bba\}$ , and  $s_0^4 = \mathbf{B}(L)$  to the language  $L$ .

Given a bitmap  $\mathbf{B} \in \{0, 1\}^{k^\ell}$ , let  $\mathbf{B}_r$  be the set of factors of  $\mathbf{B}$  of length  $k^r$ , for  $r \in [\ell]$ , in which there is at least one bit different than zero, that is,

$$\mathbf{B}_r = \{s \in \{0, 1\}^{k^r} \mid \exists i \in [k^{\ell-r} - 1] : s = s_i^r \wedge s_i^r \neq 0^{k^r}\}.$$

*Example 3.* Consider the bitmap of Example 1,  $B = 1011011100011110$ , with  $k = 2$  and  $\ell = 4$ . We have  $B_0 = \{1\}$ ,  $B_1 = \{01, 10, 11\}$ ,  $B_2 = \{0001, 0111, 1011, 1110\}$ ,  $B_3 = \{00011110, 10110111\}$ , and  $B_4 = \{B\}$ .

The size of  $B_r$  is bounded by the number of factors with size  $k^r$ . Additionally, each factor is a composition of factors from the previous set. These two conditions are formally stated in the following lemma.

**Lemma 2.** *Let  $L \subseteq \Sigma^\ell$  be a block language of words of length  $\ell > 0$  over an alphabet  $\Sigma$  of size  $k > 0$ , with a correspondent bitmap  $B$ . Then, the cardinality of  $B_r$  is bounded by:*

$$|B_r| \leq \begin{cases} 1, & \text{if } r = 0; \\ \min(k^{\ell-r}, (|B_{r-1}| + 1)^k - 1), & \text{otherwise.} \end{cases}$$

The sets  $B_r$  are related to the states of the finite automata representing the block language with bitmap  $B$ , as described in the next sections. A finite automaton for a block language is, of course, also acyclic and ranked. All final states belong to the rank  $Q_0$  and, therefore, can be merged. Additionally, if an NFA for a block language has multiple initial states, they can also be merged.

### 3.1 Minimal DFAs for Block Languages

In this section, we relate the bitmap of a block language to its minimal DFA. Given a bitmap  $B$  associated with a block language  $L \subseteq \Sigma^\ell$ , with  $\Sigma$  of size  $k$ , one can directly build a minimal DFA  $\mathcal{A}$  for  $L$ . Let  $Q = \bigcup_{r \in [\ell]} B_r$  be the set of states of  $\mathcal{A}$ , and the transition function mapping the states in  $B_r$  with the ones in  $B_{r-1}$ ,  $r \in [1, \ell]$ . We will now detail this construction.

We start by the final state  $q_f$ , that is the factor  $1 \in B_0$ , as well as the dead state corresponding to the factor 0. Then, for each rank  $r = 1, \dots, \ell$ , we consider every factor  $s \in B_r$  as a state in rank  $r$ . As stated in Lemma 1, every factor  $s$  corresponds to the bitmap of the quotient of the language  $L$  by some word  $w$ . The transitions from  $s$  are then given by the decomposition of  $s$  by  $s_0 \cdots s_{k-1}$ , where  $|s_i| = k^{\ell-r-1}$  for every  $i \in [k-1]$ . Then, for each  $i \in [k-1]$ ,  $\delta(s, \sigma_i) = s_i$ , where  $s_i \in B_{r-1}$ . Note that this construction creates one initial and one final state, since  $|B_0| = |B_r| = 1$ , if the language  $L$  is not empty.

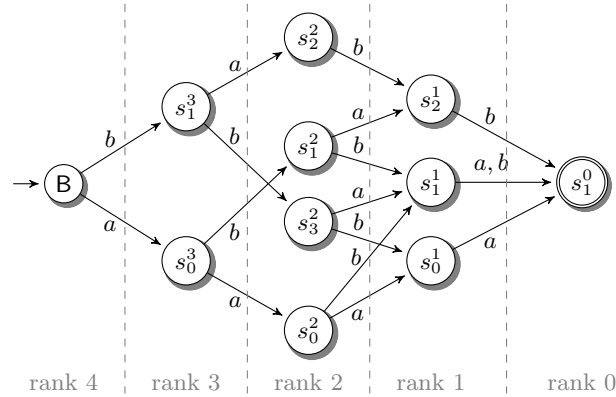
**Lemma 3.** *Let  $L \subseteq \Sigma^\ell$  be a block language,  $\ell > 0$ , with bitmap  $B$ . Then, the DFA  $\mathcal{A}$  obtained by applying the above construction to  $B$  accepts  $L$ , that is,  $\mathcal{L}(\mathcal{A}) = L$ .*

**Lemma 4.** *Let  $L \subseteq \Sigma^\ell$  be a block language,  $\ell > 0$ , with bitmap  $B$ . Then, the DFA  $\mathcal{A}$  obtained by applying the above construction to  $B$  is minimal.*

Combining the results of Lemmas 3 and 4, we obtain:

**Theorem 1.** *Given a block language  $L \subseteq \Sigma^\ell$ , the construction of a DFA from the bitmap  $B(L)$ , yields the minimal DFA for  $L$ .*

*Example 4.* Let  $L \subseteq \{a, b\}^4$  be the language of Example 1 with bitmap  $B(L) = 1011011100011110$ . The correspondent minimal DFA is depicted in Fig. 1. The final state (rank 0) labeled by  $s_1^0$  represents the bitmap factor 1 and the dead state represents  $s_0^0 = 0$  (but it is not showed, as well as all transitions from and to it). States in rank 1 correspond to 2-bit factors, in this case:  $s_0^1 = 10$ ,  $s_1^1 = 11$ , and  $s_2^1 = 01$ . We have  $\delta(s_0^1, a) = s_0^0$ ,  $\delta(s_0^1, b) = s_0^0$ , etc. States in rank 2 correspond to  $2^2 = 4$ -bit words:  $s_0^2 = 1011$ ,  $s_1^2 = 0111$ ,  $s_2^2 = 0001$  and  $s_3^2 = 1110$ . And we have, for instance,  $\delta(s_0^2, a) = s_0^1$  and  $\delta(s_0^2, b) = s_1^1$ . Similarly for ranks 3 and 4. The initial state corresponds to  $B$ .



**Fig. 1.** Minimal DFA accepting the language of Example 1, where  $s_1^0 = 1$ ,  $s_0^1 = 10$ ,  $s_1^1 = 11$ ,  $s_2^1 = 01$ ,  $s_0^2 = 1011$ ,  $s_1^2 = 0111$ ,  $s_2^2 = 0001$ ,  $s_3^2 = 1110$ ,  $s_0^3 = 10110111$ ,  $s_1^3 = 00011110$ , and  $B = 1011011100011110$ .

Given a block language  $L \subseteq \Sigma^\ell$  over an alphabet of size  $k$ , Câmpeanu and Ho [3] showed that the state complexity of a block language is at most  $\frac{k^{\ell-r}-1}{k-1} + \sum_{n=0}^{r-1} (2^{k^n} - 1)$ , where  $r = \min\{n \in [\ell] \mid k^{\ell-n} \leq 2^{k^n} - 1\}$ . In the next theorem we give an estimation of the value of  $r$ .

**Theorem 2.** *Let  $k > 1$  and  $r = \min\{n \in [\ell] \mid k^{\ell-n} \leq 2^{k^n} - 1\}$ . Then,  $r = \lceil \log_k \ell \rceil + 1 + x$ , for some  $x \in \{-1, 0, 1\}$ .*

### 3.2 Minimal NFAs for Block Languages

We now show that the bitmap of a block language  $L$  can be used to obtain a minimal NFA for  $L$ . However in this case the problem is NP-complete. Given a bitmap  $B$  associated with a block language  $L \subseteq \Sigma^\ell$ , one can build a minimal NFA similarly to the previous construction for minimal DFAs, by iteratively finding the minimal number of states required at each rank. The main difference

with the deterministic case is that the quotients of the language, corresponding to factors from the bitmap, are represented by a set of states, instead of a single one.

Let us first define what is a *cover* with respect to a bitmap factor. Consider the language  $L_r$  as the set of bitmaps in  $\{0, 1\}^{k^r}$  which can be split into  $k$  subwords of length  $k^{r-1}$ , where there is at least one subword which is in  $B_{r-1}$ . That is,  $L_r$  is the language  $(B_{r-1} + 0^{k^{r-1}})^k \setminus 0^{k^r}$ . We say that  $C_r \subseteq L_r$  is a cover for a word  $w \in \{0, 1\}^{k^r}$  (or, alternatively,  $w$  is covered by  $C_r$ ) if there is a subset  $\{c_0, \dots, c_{n-1}\} \subseteq C_r$  with  $n \leq |C_r|$ , such that  $\bigvee_{i \in [n-1]} c_i = w$ . For example, the sequence 1011 is covered by the set  $\{0011, 1010\}$  because  $1011 = 0011 \vee 1010$ . Moreover,  $C_r \subseteq L_r$  is a cover for the set  $B_r$  if all the words in  $B_r$  are covered by  $C_r$ . For instance, it can be easily noticed that  $B_r$  covers itself.

We now show how to obtain a nondeterministic finite automaton for  $L$  starting from  $B$ . The construction starts as before, where the final state  $q_f$  and the dead state, are the bitmap factors 1 and 0, respectively. Additionally, we define the function  $\rho : \{0, 1\}^* \rightarrow 2^{\{0,1\}^*}$  which will map each bitmap into the set that covers it, initially set to  $\rho(0) = \{0\}$  and  $\rho(1) = \{1\}$ .

Next, for every rank  $r = 1, \dots, \ell$ , we consider  $C_r$  as is a minimal cover for  $B_r$ , and we set, for every  $s \in B_r$ ,  $\rho(s) = \{c_0, \dots, c_{n-1}\} \subseteq C_r$ , such that  $\rho(s)$  covers  $s$ . Then, take  $C_r$  as the set of states in rank  $r$  and, as in the deterministic construction, with the transitions from state  $s$  according to  $c$ .

More precisely, let  $\sigma_i \in \Sigma$  and  $c_i$  be the  $i$ -th subword of  $c$  of length  $k^{r-1}$ , for  $i \in [k-1]$ . Let  $c \in C_r$ , then,  $\delta(c, \sigma_i) = \{s' \in C_{r-1} \mid s' \in \rho(c_i)\}$ .

**Lemma 5.** *Given a block language  $L \subseteq \Sigma^\ell$ , for some  $\ell > 0$ , with bitmap  $B$ , then the NFA  $\mathcal{A}$  given by the above construction applied to  $B$  accepts  $L$ , that is,  $\mathcal{L}(\mathcal{A}) = L$ .*

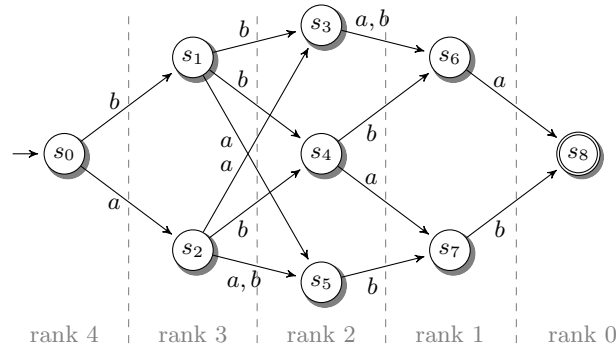
**Lemma 6.** *Given a block language  $L \subseteq \Sigma^\ell$ , for some  $\ell > 0$ , with bitmap  $B$ , then the NFA  $\mathcal{A}$  given by the above construction applied to  $B$  is minimal.*

We obtain the following result from Lemmas 5 and 6.

**Theorem 3.** *The construction from a bitmap, with respect to a block language  $L \subseteq \Sigma^\ell$ , to a NFA returns a DFA  $\mathcal{A}$  such that  $\mathcal{L}(\mathcal{A}) = L$ , furthermore,  $\mathcal{A}$  is minimal.*

*Example 5.* Let  $L \subseteq \{a, b\}^4$  be the language of Example 1 with bitmap  $B = 1011011100011110$ . A correspondent NFA  $\mathcal{N}$  is depicted in Fig. 2. We have that  $B_1 = \{01, 11, 10\}$  but  $C_1 = \{01, 10\}$  is a cover so only two states are needed in rank 1. Then,  $B_2 = \{1011, 0111, 0001, 1110\}$ , and take  $C_2 = \{1010, 0110, 0001\}$ . We have that

- $1011 = 0001 \vee 1010$
- $0111 = 0110 \vee 0001$
- $1110 = 1010 \vee 0110$



**Fig. 2.** A minimal NFA accepting the language of Example 1, where  $s_8 = 1$ ,  $s_7 = 01$ ,  $s_6 = 10$ ,  $s_5 = 0001$ ,  $s_4 = 0110$ ,  $s_3 = 1010$ ,  $s_2 = 10110111$ ,  $s_1 = 00011110$ , and  $s_0 = 1011011100011110$ .

In rank 3 we need at least two states and rank 4 has only the initial state.

This problem is NP-COMPLETE, as the following theorem formalises.

**Theorem 4.** *Given a set of factors  $B_r$ , each with size  $k^r$  for some integers  $k, r > 0$ , the problem of finding the minimal cover  $C_r$  such that each element from  $B_r$  can be obtained by the disjunction of a subset from  $C_r$  is NP-COMPLETE.*

One can use an SMT-solver [11] to find a cover of size  $n$  of a set, wherein every factor of size  $k^r$  can be represented as an equally sized bit vector, and then employ binary search on the solution to determine the minimal one.

### 3.3 Maximal Size of Minimal NFAs for Block Languages

The width of each rank on the NFA given by the construction described above must be bounded by the width of the same rank on the DFA. Also, to cover factors of length  $k^r$ , we show that a cover of  $k^r$  elements is sufficient. These bounds are formally stated on the following lemma.

**Lemma 7.** *Let  $\mathcal{A} = \langle Q, \Sigma, \delta, q_0, \{q_f\} \rangle$  be a NFA for a block language  $L \subseteq \Sigma^\ell$  over  $\Sigma = \{\sigma_0, \dots, \sigma_{k-1}\}$  of size  $k$ , given by the construction of NFAs from bitmaps. The set of states  $Q$  can be partitioned into  $\ell + 1$  sets  $Q_0 \cup \dots \cup Q_\ell$ , such that every state  $q \in Q_r$  has  $\text{rank}(q) = r$ , for  $r \in [\ell]$ . Then, the cardinality of the set of states at rank  $r$ , namely  $|Q_r|$ , is bounded by:*

$$|Q_r| \leq \min(k^{\ell-r}, k^r), \quad \text{for all } r \in [\ell].$$

Lemma 7 allows us to determine the exact maximal size of a minimal NFA for a block language, stated in the following theorem.



**Theorem 5.** *The maximal size of a minimal NFA for a block language  $L \subseteq \Sigma^\ell$ , with  $\ell > 0$  and  $|\Sigma| = k$ , is  $nsc(L) \leq 2 \cdot \frac{k^{\frac{\ell}{2}} - 1}{k - 1} + k^{\frac{\ell}{2}}$  if  $\ell$  is even, and  $nsc(L) \leq 2 \cdot \frac{k^{\lceil \frac{\ell}{2} \rceil} - 1}{k - 1}$ , otherwise.*

## 4 Operational Complexity

In this section we consider several operations on block languages using its bitmap representations and study the state complexity of those operations. Formally, the *operational state complexity* is the worst-case state complexity of a language resulting from the operation, considered as a function of the state complexities of the operands. We will focus on operations that are closed for block languages, i.e., the resulting language is also a block language.

In general, the upper bounds of operational state complexity known for finite languages apply for block languages. Let  $L_1$  and  $L_2$  be two block languages with  $sc(L_1) = n$  and  $sc(L_2) = m$ . For the concatenation the upper bound for  $sc(L_1L_2)$  coincide with the bound for finite languages when the minimal DFA of  $L_1$  has only one final state (a part from the initial state), that is  $m + n - 1$ . For union and intersection when  $L_1, L_2 \subseteq \Sigma^\ell$  the bounds are  $O(mn)$  but we give a more precise value. For reversal we present a witness for which the upper bound is exponential. Finally we consider the operational state complexity of the set difference  $\Sigma^\ell \setminus L_1$ , for  $L_1 \subseteq \Sigma^\ell$ .

*Reversal.* Given a bitmap  $\mathbf{B}$  of a block language  $L \subseteq \Sigma^\ell$  and  $|\Sigma| = k$ , in the following we compute the bitmap for the reversal language  $L^R$ .

Recall that the perfect shuffle of length 1,  $\sqcup_1$ , of two words  $u = u_0 \cdots u_{n-1}$  and  $v = v_0 \cdots v_{n-1}$  of the same length is given by  $u \sqcup_1 v = u_0v_0 \cdots u_{n-1}v_{n-1}$ . If  $j$  is a divisor of  $n$ , the perfect shuffle of length  $j$  of  $u$  and  $v$  is the perfect shuffle of blocks of length  $j$ , i.e.,  $u \sqcup_j v = u_0 \cdots u_{j-1}v_0 \cdots v_{j-1} \cdots u_{n-j} \cdots u_{n-1}v_{n-j} \cdots v_{n-1}$ . Finally, if  $|u| = |v|$ , we denote  $u \sqcup_j v$  by  $\sqcup_j^2(uv)$ . This can be generalised for any number  $m \geq 2$  of words  $w_0, \dots, w_{m-1}$  of the same length by considering the perfect shuffle of blocks of length  $j$  taken from each of the  $w_i$  words, that is,  $\sqcup_j^m(w_0 \cdots w_{k-1})$ . For  $j = 1$  and  $w_i = (w_{i,0} \cdots w_{i,n-1})$ ,  $i \in [k - 1]$ , one gets

$$\sqcup_1^k(w_0 \cdots w_{k-1}) = w_{0,0} \cdots w_{k-1,0} w_{0,1} \cdots w_{k-1,1} \cdots w_{0,n-1} \cdots w_{k-1,n-1}.$$

Let  $\mathbf{R}_0 = \mathbf{B}$  and  $\mathbf{R}_i = \sqcup_i^k(\mathbf{R}_{i-1})$ , for  $i \in [1, \ell - 1]$ .

**Lemma 8.** *Let  $L \subseteq \Sigma^\ell$ , for some  $\ell > 0$ . The bitmap for the reversal of  $L$ , namely  $\mathbf{B}(L^R)$ , can be obtained with  $\mathbf{R}_{\ell-1}$ .*

*Example 6.* Consider  $\Sigma = \{a, b\}$  and  $\ell = 3$ . Any bitmap is of the form  $\mathbf{B} = b_0b_1b_2b_3b_4b_5b_6b_7$ . We have  $\mathbf{R}_1 = b_0b_4b_1b_5b_2b_6b_3b_7$  and  $\mathbf{R}_2 = b_0b_4b_2b_6b_1b_5b_3b_7$ , which is the bitmap of the reversal language.

We now study the state complexity of the reversal of a block language by giving an witness family parametrized by the block length  $\ell$  for which the reversal has a state complexity that is exponential in the state complexity of the witness.

**Lemma 9.** *Consider the language  $L \subseteq \{a, b\}^\ell$ , for some  $\ell > 0$ , with bitmap*

$$\mathbb{B}(L) = \prod_{i=1}^{n_r} \text{pad}(i_{[2]}, 2^r)^R \cdot 0^{2^r \cdot (2^{\ell-r} - n_r)}$$

where  $i_{[2]}$  denotes the binary representation of  $i$ ,  $\text{pad}(n, m)$  the function which adds leading zeros to the string  $n$  until its length is  $m$ ,  $n_r$  the width of a minimal DFA with maximal size for a block language in  $\Sigma^\ell$  and  $r \in [\ell]$  its respective rank.

Then, the minimal DFA which accepts the language  $L$  has maximal size.

The language  $L$  represented by the bitmap of the previous lemma can be formally described as

$$L = \{w_1 w_2 \mid w_1 \in \Sigma^{\ell-r}, w_2 \in \Sigma^r, i = \text{ind}(w_1), j = \text{ind}(w_2), (i+1) \wedge 2^j \neq 0\}$$

where  $\text{ind}(w)$  represents the index of  $w$  in  $\Sigma^{|w|}$ , in a lexicographic order. Informally, the words from the language can be split into two subwords,  $w_1$  and  $w_2$ , with lengths  $\ell - r$  and  $r$ , respectively. Moreover, their indexes,  $i = \text{ind}(w_1)$  and  $j = \text{ind}(w_2)$ , must imply that the  $j$ -th bit of the binary representation of  $(i+1)$  is set to 1.

*Example 7.* The rank widths of the minimal DFA with maximal size for a language of block length  $\ell = 5$  and and alphabet size  $k = 2$  is  $\{1, 3, 8, 4, 2, 1\}$ . The width of the DFA is  $n_r = 8$  obtained at the rank  $r = 2$ . The bitmap of the language  $L$ , according to Lemma 9, is 10000100110000101010011011100001 and

$$L = \{aaaaa, aabab, abaaa, abaab, abbba, baaaa, baaba, babab, babba, bbaaa, bbaab, bbaba, bbbbb\}$$

Let  $w_1 = baa$  where  $i = \text{ind}(w_1) = 4$  and the binary representation of  $i+1$  is 101. With  $j$  equals 0 and 2, we have that  $(i+1) \wedge 2^j \neq 0$ , which corresponds to the words  $aa$  and  $ba$ . Therefore  $baaaa, baaba \in L$ .

**Lemma 10.** *Let  $r, n_r, \ell$ , and  $L$  be the same as in Lemma 9. A DFA whose width is limited by the width of rank  $\ell - r - 1$  is sufficient to accept the reverse of the language  $L$ .*

**Lemma 11.** *Let  $r, n_r$ , and  $\ell$  be the same from Lemma 9. Let  $\mathcal{A}$  and  $\mathcal{B}$  be two DFAs for block languages in  $\{a, b\}^\ell$  such that  $\text{width}(\mathcal{A}) = r$  and  $\text{width}(\mathcal{B}) = \ell - r - 1$ . Then the number of states of  $\mathcal{A}$  is exponential in the number of states of  $\mathcal{B}$ .*

**Theorem 6.** *Let  $L \subseteq \Sigma^\ell$ ,  $\ell > 0$ , such that  $sc(L) = n$ . Then,  $sc(L^R)$  is  $2^{O(n)}$ .*

*Intersection.* Given two languages  $L_1, L_2 \subseteq \Sigma^\ell$ , for some  $\ell > 0$  and  $|\Sigma| = k$ , and their respective bitmaps  $\mathbf{B}(L_1), \mathbf{B}(L_2)$ , the bitmap of the intersection of both languages is given by  $\mathbf{B}(L_1) \wedge \mathbf{B}(L_2)$ .

**Theorem 7.** *Let  $L_1, L_2 \subseteq \Sigma^\ell$  be two block languages with  $\ell > 0$ . Also, let  $\mathcal{A} = \langle Q, \Sigma, \delta_1, q_0, \{q_f\} \rangle$  and  $\mathcal{B} = \langle P, \Sigma, \delta_2, p_0, \{p_f\} \rangle$ , two minimal DFAs such that  $\mathcal{L}(\mathcal{A}) = L_1$  and  $\mathcal{L}(\mathcal{B}) = L_2$ . Moreover,  $\mathcal{A}$  and  $\mathcal{B}$  have rank widths  $\{n_0, n_1, \dots, n_\ell\}$  and  $\{m_0, m_1, \dots, m_\ell\}$ , respectively, where  $n_0 = m_0 = n_\ell = m_\ell = 1$ ,  $\sum_{r=0}^{\ell} n_r = |Q|$  and  $\sum_{r=0}^{\ell} m_r = |P|$ . Then,  $sc(L_1 \cap L_2) \leq \sum_{r=0}^{\ell} n_r m_r$ , and the bound is tight if the size of the alphabet can depend on the sizes of  $\mathcal{A}$  and  $\mathcal{B}$ .*

*Union.* Given two languages  $L_1, L_2 \subseteq \Sigma^\ell$ , for some  $\ell > 0$  and  $|\Sigma| = k$ , and their respective bitmaps  $\mathbf{B}(L_1), \mathbf{B}(L_2)$ , the bitmap of the union of both languages is given by  $\mathbf{B}(L_1) \vee \mathbf{B}(L_2)$ .

**Theorem 8.** *Let  $L_1, L_2 \subseteq \Sigma^\ell$  be two block languages with  $\ell > 0$ . Also, let  $\mathcal{A} = \langle Q, \Sigma, \delta_1, q_0, \{q_f\} \rangle$  and  $\mathcal{B} = \langle P, \Sigma, \delta_2, p_0, \{p_f\} \rangle$ , two minimal DFAs such that  $\mathcal{L}(\mathcal{A}) = L_1$  and  $\mathcal{L}(\mathcal{B}) = L_2$ . Moreover,  $\mathcal{A}$  and  $\mathcal{B}$  have rank widths  $\{n_0, n_1, \dots, n_\ell\}$  and  $\{m_0, m_1, \dots, m_\ell\}$ , respectively, where  $n_0 = m_0 = n_\ell = m_\ell = 1$ ,  $\sum_{r=0}^{\ell} n_r = |Q|$  and  $\sum_{r=0}^{\ell} m_r = |P|$ . Then,  $sc(L_1 \cup L_2) \leq 2 + \sum_{r=1}^{\ell-1} (n_r m_r + n_r + m_r)$ , and the bound is tight if the size of the alphabet can depend on the sizes of  $\mathcal{A}$  and  $\mathcal{B}$ .*

*Concatenation.* Consider two languages  $L_1 \subseteq \Sigma^{\ell_1}$  and  $L_2 \subseteq \Sigma^{\ell_2}$ , for some  $\ell_1, \ell_2 > 0$ , with bitmaps  $\mathbf{B}(L_1)$  and  $\mathbf{B}(L_2)$ , respectively. The bitmap for the language  $L_1 L_2$  can be constructed as follows: For  $i = 0, \dots, k^{\ell_1}$ , if  $b_i = 1$  then append  $\mathbf{B}(L_2)$ , otherwise append  $0^{k^{\ell_2}}$ . Then, we have the following exact state complexity for concatenation:

**Theorem 9.** *Let  $L_1 \subseteq \Sigma^{\ell_1}$  and  $L_2 \subseteq \Sigma^{\ell_2}$ , for some  $\ell_1, \ell_2 > 0$  over an alphabet of size  $k$ . Then,  $sc(L_1 L_2) = sc(L_1) + sc(L_2) - 1$ .*

*Block Complement.* Consider a language  $L \subseteq \Sigma^\ell$ , for some  $\ell > 0$  and alphabet of size  $k > 0$ . Let  $\mathbf{B}$  be its bitmap. Then, the bitmap of the language  $\Sigma^\ell \setminus L$ , namely  $\overline{\mathbf{B}}$ , is given by flipping every bit of  $\mathbf{B}$ .

**Theorem 10.** *Let  $L \subseteq \Sigma^\ell$ , with  $\ell > 0$ , be a block language with  $|\Sigma| = k$ , such that  $sc(L) = n$ . Then,  $n - (\ell - 1) \leq sc(\Sigma^\ell \setminus L) \leq n + (\ell - 1)$ .*

## References

1. Almeida, M., Moreira, N., Reis, R.: Exact generation of minimal acyclic deterministic finite automata. *Int. J. Found. Comput. S.* **19**(4), 751–765 (2008). <https://doi.org/10.1142/S0129054108005930>
2. Almeida, M., Moreira, N., Reis, R.: Finite automata minimization algorithms. In: Wang, J. (ed.) *Handbook of Finite State Based Models and Applications*, pp. 145–170. CRC Press (2012), <http://www.crcpress.com/product/isbn/9781439846186>

3. Câmpeanu, C., Ho, W.H.: The maximum state complexity for finite languages. *J. Autom. Lang. Comb.* **9**(2-3), 189–202 (2004)
4. Câmpeanu, C., Culik II, K., Salomaa, K., Yu, S.: State complexity of basic operations on finite languages. In: Boldt, O., Jürgensen, H. (eds.) 4th WIA'99. LNCS, vol. 2214, pp. 60–70. Springer-Verlag (2001)
5. Elvey Price, A., Fang, W., Wallner, M.: Compacted binary trees admit a stretched exponential. *J. Comb. Theory, Ser. A* **177**, 105306 (2021). <https://doi.org/10.1016/J.JCTA.2020.105306>
6. Gao, Y., Moreira, N., Reis, R., Yu, S.: A survey on operational state complexity. *Journal of Automata, Languages and Combinatorics* **21**(4), 251–310 (2017)
7. Karhumäki, J., Kari, J.: Finite automata, image manipulation, and automatic real functions. In: Pin, J. (ed.) *Handbook of Automata Theory*, pp. 1105–1143. European Mathematical Society (2021). <https://doi.org/10.4171/AUTOMATA-2/8>
8. Karhumäki, J., Okhotin, A.: On the determinization blowup for finite automata recognizing equal-length languages. In: C. S. Calude, R.F., Iwama, K. (eds.) *Computing with New Resources - Essays Dedicated to Jozef Gruska*. LNCS, vol. 8808, pp. 71–82. Springer (2014). [https://doi.org/10.1007/978-3-319-13350-8\\_6](https://doi.org/10.1007/978-3-319-13350-8_6)
9. Kjos-Hanssen, B., Liu, L.: The number of languages with maximum state complexity. In: Gopal, T.V., Watada, J. (eds.) 15th TAMC. LNCS, vol. 11436, pp. 394–409. Springer (2019). [https://doi.org/10.1007/978-3-030-14812-6\\_24](https://doi.org/10.1007/978-3-030-14812-6_24)
10. Konstantinidis, S., Moreira, N., Reis, R.: Randomized generation of error control codes with automata and transducers. *RAIRO* **52**, 169–184 (2018)
11. Kroening, D., Strichman, O.: *Decision Procedures: An Algorithmic Point of View*. Springer (2016). <https://doi.org/10.1007/978-3-662-50497-0>
12. Meyer, A.R., Fischer, M.J.: Economy of description by automata, grammars, and formal systems. In: 12th Annual Symposium on Switching and Automata Theory. pp. 188–191. IEEE, Los Alamitos (1971)
13. Priez, J.B.: Enumeration of minimal acyclic automata via generalized parking functions. In: 27th FPSAC. DMTCS, vol. 2471 (2015), <https://doi.org/10.46298/dmtcs.2471>
14. Revuz, D.: Minimisation of acyclic deterministic automata in linear time. *Theoret. Comput. Sci.* **92**(1), 181–189 (1992)
15. Stockmeyer, L., Meyer, A.R.: Word problems requiring exponential time: Preliminary report. In: 5th Annual ACM Symposium on Theory of Computing. pp. 1–9. ACM (1973)
16. Yu, S., Zhuang, Q., Salomaa, K.: The state complexities of some basic operations on regular languages. *Theor. Comput. Sci.* **125**(2), 315–328 (1994)