

Unrestricted State Complexity of Binary Operations on Regular Languages

Janusz Brzozowski

► **To cite this version:**

Janusz Brzozowski. Unrestricted State Complexity of Binary Operations on Regular Languages. 18th International Workshop on Descriptive Complexity of Formal Systems (DCFS), Jul 2016, Bucharest, Romania. pp.60-72, 10.1007/978-3-319-41114-9_5. hal-01633951

HAL Id: hal-01633951

<https://hal.inria.fr/hal-01633951>

Submitted on 13 Nov 2017

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Unrestricted State Complexity of Binary Operations on Regular Languages^{*}

Janusz Brzozowski

David R. Cheriton School of Computer Science, University of Waterloo,
Waterloo, ON, Canada N2L 3G1
{brzozo@uwaterloo.ca}

Abstract. I study the state complexity of binary operations on regular languages over different alphabets. It is well known that if L'_m and L_n are languages restricted to be over the same alphabet, with m and n quotients, respectively, the state complexity of any binary boolean operation on L'_m and L_n is mn , and that of the product (concatenation) is $(m-1)2^n + 2^{n-1}$. In contrast to this, I show that if L'_m and L_n are over their own different alphabets, the state complexity of union and symmetric difference is $mn + m + n + 1$, that of intersection is $mn + 1$, that of difference is $mn + m + 1$, and that of the product is $m2^n + 2^{n-1}$.

Keywords: boolean operation, concatenation, different alphabets, most complex languages, product, quotient complexity, regular language, state complexity, stream, unrestricted complexity

1 Motivation

Formal definitions are postponed until Section 2.

The first paper on state complexity was published by A. N. Maslov [9] in 1970, but this work was unknown in the West for many years. Maslov wrote:

An important measure of the complexity of [sets of words representable in finite automata] is the number of states in the minimal representing automaton. ... if $T(A) \cup T(B)$ are representable in automata A and B with m and n states respectively ..., then:

- 1. $T(A) \cup T(B)$ is representable in an automaton with $m \cdot n$ states;*
- 2. $T(A).T(B)$ is representable in an automaton with $(m-2)2^n + 2^{n-1}$ states.*

The second paper on state complexity was published by S. Yu, Q. Zhuang and K. Salomaa [11] in 1994. Here the authors wrote:

- 1. ... for any pair of complete m -state DFA A and n -state DFA B defined on the same alphabet Σ , there exists a DFA with at most $m2^n - 2^{n-1}$ states which accepts $L(A)L(B)$.*

^{*} This work was supported by the Natural Sciences and Engineering Research Council of Canada grant No. OGP0000871.

2. ... $m \cdot n$ states are ... sufficient for a DFA to accept the intersection (union) of an m -state DFA language and an n -state DFA language.

Here DFA stands for *deterministic finite automaton*, and *complete* means that there is a transition from every state under every input letter.

I will show that statements 1 and 2 of Maslov are incorrect without the restriction that the languages are over the same alphabet. In [11] the first statement includes that restriction, but the second omits it (presumably it's implied).

The same-alphabet restriction is unnecessary: There is no reason why we should not be able to find, for example, the union of languages $L' = \{a, b\}^*b$ and $L = \{a, c\}^*c$ accepted by the minimal complete two-state automata \mathcal{D}'_2 and \mathcal{D}_2 of Figure 1, where an incoming arrow denotes the initial state and a double circle represents a final state.



Fig. 1. Two minimal complete DFAs \mathcal{D}'_2 and \mathcal{D}_2 .

The union of L' and L is a language over three letters. To find the DFA for $L' \cup L$, we view \mathcal{D}'_2 and \mathcal{D}_2 as incomplete DFA's, the first missing all transitions under c , and the second under b . After adding the missing transitions we obtain DFAs \mathcal{D}'_3 and \mathcal{D}_3 . Now we can proceed as is usually done in the same-alphabet approach, and take the direct product of \mathcal{D}'_3 and \mathcal{D}_3 to find $L' \cup L$. Here it turns out that six states are necessary to represent $L' \cup L$, but the state complexity of union is actually $(m + 1)(n + 1)$.



Fig. 2. DFAs \mathcal{D}'_3 and \mathcal{D}_3 over three letters.

In general, when calculating the result of a binary operation on regular languages with different alphabets, we deal with special incomplete DFAs that are

only missing some letters and all the transitions caused by these letters. The complexity of incomplete DFAs has been studied previously by Gao, K. Salomaa, and Yu [6] and by Maia, Moreira and Reis [8]. However, the objects studied there are *arbitrary* incomplete DFAs, whereas we are interested only in *complete DFAs with some missing letters*. Secondly, we study *state* complexity, whereas the above-mentioned papers deal mainly with *transition* complexity. Nevertheless, there is some overlap. It was shown in [6, Corollary 3.2] that the incomplete state complexity of union is less than or equal to $mn + m + n$, and that this bound is tight in some special cases. In [8, Theorem 2], witnesses that work in all cases were found. These complexities correspond to my result for union in Theorem 1. Also in [8, Theorem 5], the incomplete state complexity of product is shown to be $m2^n + 2^{n-1} - 1$, and this corresponds to my result for product in Theorem 2.

In this paper I remove the restriction of equal alphabets of the two operands. I prove that the complexity of union and symmetric difference is $mn + m + n + 1$, that of intersection is $mn + 1$, that of difference is $mn + m - 1$, and that of the product is $m2^n + 2^{n-1}$, if each language's own alphabet is used. I exhibit a new most complex regular language that meets the complexity bounds for boolean operations, product, star, and reversal, has a maximal syntactic semigroup and most complex atoms. All the witnesses used here are derived from that one most complex language.

2 Terminology and Notation

A basic complexity measure of a regular language L over an alphabet Σ is the number n of distinct (left) quotients of L , where a (*left*) *quotient* of L by a word $w \in \Sigma^*$ is $w^{-1}L = \{x \mid wx \in L\}$. The number of quotients of L is its *quotient complexity* [2], $\kappa(L)$. A concept equivalent to quotient complexity is the *state complexity* [11] of L , which is the number of states in a complete minimal deterministic finite automaton (DFA) recognizing L . Since we do not use any other measures of complexity in this paper (with the exception of one mention of time and space complexity in the next paragraph), we refer to quotient/state complexity simply as *complexity*.

Let $L'_m \subseteq \Sigma'^*$ and $L_n \subseteq \Sigma^*$ be regular languages of complexities m and n , respectively. The *complexity of a binary operation* \circ on L'_m and L_n is the maximal value of $\kappa(L'_m \circ L_n)$ as a function $f(m, n)$, as L'_m and L_n range over all regular languages of complexity m and n , respectively. The complexity of an operation gives a worst-case lower bound on the time and space complexity of the operation. For this reason it has been studied extensively; see [2, 3, 10, 11] for additional references.

A *deterministic finite automaton (DFA)* is a quintuple $\mathcal{D} = (Q, \Sigma, \delta, q_0, F)$, where Q is a finite non-empty set of *states*, Σ is a finite non-empty *alphabet*, $\delta: Q \times \Sigma \rightarrow Q$ is the *transition function*, $q_0 \in Q$ is the *initial state*, and $F \subseteq Q$ is the set of *final states*. We extend δ to a function $\delta: Q \times \Sigma^* \rightarrow Q$ as usual. A DFA \mathcal{D} *accepts* a word $w \in \Sigma^*$ if $\delta(q_0, w) \in F$. The language accepted by \mathcal{D} is

denoted by $L(\mathcal{D})$. If q is a state of \mathcal{D} , then the language L^q of q is the language accepted by the DFA $(Q, \Sigma, \delta, q, F)$. A state is *empty* (or *dead* or a *sink state*) if its language is empty. Two states p and q of \mathcal{D} are *equivalent* if $L^p = L^q$. A state q is *reachable* if there exists $w \in \Sigma^*$ such that $\delta(q_0, w) = q$. A DFA is *minimal* if all of its states are reachable and no two states are equivalent. Usually DFAs are used to establish upper bounds on the complexity of operations, and also as witnesses that meet these bounds.

If $\delta(q, a) = p$ for a state $q \in Q$ and a letter $a \in \Sigma$, we say there is a *transition* under a from q to p in \mathcal{D} . The DFAs defined above are *complete* in the sense that there is *exactly one* transition for each state $q \in Q$ and each letter $a \in \Sigma$. If there is *at most one transition* for each state of Q and letter of Σ , the automaton is an *incomplete* DFA.

A *nondeterministic finite automaton (NFA)* is a 5-tuple $\mathcal{D} = (Q, \Sigma, \delta, I, F)$, where Q, Σ and F are defined as in a DFA, $\delta: Q \times \Sigma \rightarrow 2^Q$ is the *transition function*, and $I \subseteq Q$ is the *set of initial states*. An ε -NFA is an NFA in which transitions under the empty word ε are also permitted.

To simplify the notation, without loss of generality we use $Q_n = \{0, \dots, n-1\}$ as the set of states of every DFA with n states. A *transformation* of Q_n is a mapping $t: Q_n \rightarrow Q_n$. The *image* of $q \in Q_n$ under t is denoted by qt . For $k \geq 2$, a transformation (permutation) t of a set $P = \{q_0, q_1, \dots, q_{k-1}\} \subseteq Q$ is a *k-cycle* if $q_0t = q_1, q_1t = q_2, \dots, q_{k-2}t = q_{k-1}, q_{k-1}t = q_0$. This *k-cycle* is denoted by $(q_0, q_1, \dots, q_{k-1})$, and acts as the identity on the states in $Q_n \setminus P$. A 2-cycle (q_0, q_1) is called a *transposition*. A transformation that changes only one state p to a state $q \neq p$ and acts as the identity for the other states is denoted by $(p \rightarrow q)$. The identity transformation is denoted by $\mathbf{1}$.

In any DFA, each $a \in \Sigma$ induces a transformation δ_a of the set Q_n defined by $q\delta_a = \delta(q, a)$; we denote this by $a: \delta_a$. For example, when defining the transition function of a DFA, we write $a: (0, 1)$ to mean that $\delta(q, a) = q(0, 1)$, where the transformation $(0, 1)$ acts on state q as follows: if q is 0 it maps it to 1, if q is 1 it maps it to 0, and it acts as the identity on the remaining states.

By a slight abuse of notation we use the letter a to denote the transformation it induces; thus we write qa instead of $q\delta_a$. We extend the notation to sets of states: if $P \subseteq Q_n$, then $Pa = \{pa \mid p \in P\}$. We also find it convenient to write $P \xrightarrow{a} Pa$ to indicate that the image of P under a is Pa . If s, t are transformations of Q , their composition is denoted by $s*t$ and defined by $q(s*t) = (qs)t$; the $*$ is usually omitted. Let \mathcal{T}_{Q_n} be the set of all n^n transformations of Q_n ; then \mathcal{T}_{Q_n} is a monoid under composition.

A sequence $(L_n, n \geq k) = (L_k, L_{k+1}, \dots)$, of regular languages is called a *stream*; here k is usually some small integer, and the languages in the stream usually have the same form and differ only in the parameter n . For example, $(\{a, b\}^* a^n \{a, b\}^* \mid n \geq 2)$ is a stream. To find the complexity of a binary operation \circ we need to find an upper bound on this complexity and two streams $(L'_m, m \geq h)$ and $(L_n, n \geq k)$ of languages meeting this bound. In general, the two streams are different, but there are many examples where L'_n “differs only slightly” from L_n ; such a language L'_n is called a *dialect* [3] of L_n .

Let $\Sigma = \{a_1, \dots, a_k\}$ be an alphabet; we assume that its elements are ordered as shown. Let π be a *partial permutation* of Σ , that is, a partial function $\pi: \Sigma \rightarrow \Gamma$ where $\Gamma \subseteq \Sigma$, for which there exists $\Delta \subseteq \Sigma$ such that π is bijective when restricted to Δ and undefined on $\Sigma \setminus \Delta$. We denote undefined values of π by “–”, that is, we write $\pi(a) = -$, if π is undefined at a .

If $L \subseteq \Sigma^*$, we denote it by $L(a_1, \dots, a_k)$ to stress its dependence on Σ . If π is a partial permutation, let $s_\pi(L(a_1, \dots, a_k))$ be the language obtained from $L(a_1, \dots, a_k)$ by the substitution s_π defined as follows: for $a \in \Sigma$, $a \mapsto \{\pi(a)\}$ if $\pi(a)$ is defined, and $a \mapsto \emptyset$ otherwise. The *permutational dialect*, or simply *dialect*, of $L(a_1, \dots, a_k)$ defined by π is the language $L(\pi(a_1), \dots, \pi(a_k)) = s_\pi(L(a_1, \dots, a_k))$.

Similarly, let $\mathcal{D} = (Q, \Sigma, \delta, q_0, F)$ be a DFA; we denote it by $\mathcal{D}(a_1, \dots, a_k)$ to stress its dependence on Σ . If π is a partial permutation, then the *permutational dialect*, or simply *dialect*, $\mathcal{D}(\pi(a_1), \dots, \pi(a_k))$ of $\mathcal{D}(a_1, \dots, a_k)$ is obtained by changing the alphabet of \mathcal{D} from Σ to $\pi(\Sigma)$, and modifying δ so that in the modified DFA $\pi(a_i)$ induces the transformation induced by a_i in the original DFA. One verifies that if the language $L(a_1, \dots, a_k)$ is accepted by DFA $\mathcal{D}(a_1, \dots, a_k)$, then $L(\pi(a_1), \dots, \pi(a_k))$ is accepted by $\mathcal{D}(\pi(a_1), \dots, \pi(a_k))$.

If the letters for which π is undefined are at the end of the alphabet Σ , then they are omitted. For example, if $\Sigma = \{a, b, c, d\}$ and $\pi(a) = b$, $\pi(b) = a$, and $\pi(c) = \pi(d) = -$, then we write $L_n(b, a)$ for $L_n(b, a, -, -)$, etc.

3 Boolean Operations

A binary boolean operation is *proper* if it is not a constant and does not depend on only one variable. We study the complexities of four proper boolean operations only: union (\cup), symmetric difference (\oplus), difference (\setminus), and intersection (\cap); the complexity of any other proper operation can be deduced from these four. For example, $\kappa(\overline{L'} \cup L) = \kappa(\overline{\overline{L'} \cup L}) = \kappa(L' \cap \overline{L}) = \kappa(L' \setminus L)$, where we have used the well-known fact that $\kappa(\overline{L}) = \kappa(L)$, for any L .

The DFA of Definition 1 is required for the next theorem; this DFA is the 4-input “universal witness” called $\mathcal{U}_n(a, b, c, d)$ in [3].

Definition 1. For $n \geq 3$, let $\mathcal{D}_n = \mathcal{D}_n(a, b, c, d) = (Q_n, \Sigma, \delta_n, 0, \{n-1\})$, where $\Sigma = \{a, b, c, d\}$, and δ_n is defined by the transformations $a: (0, \dots, n-1)$, $b: (0, 1)$, $c: (n-1 \rightarrow 0)$, and $d: \mathbf{1}$. Let $L_n = L_n(a, b, c, d)$ be the language accepted by \mathcal{D}_n . The structure of $\mathcal{D}_n(a, b, c, d)$ is shown in Figure 3.

Theorem 1. For $m, n \geq 3$, let L'_m (respectively, L_n) be a regular language with m (respectively, n) quotients over an alphabet Σ' (respectively, Σ). Then $\kappa(L'_m \cup L_n) = \kappa(L'_m \oplus L_n) = mn + m + n + 1$, $\kappa(L'_m \setminus L_n) = mn + m + 1$, $\kappa(L'_m \cap L_n) = mn + 1$.

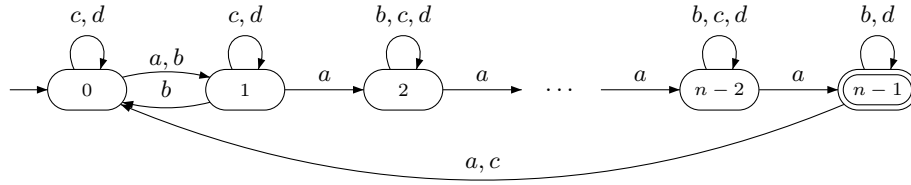


Fig. 3. DFA of Definition 1.

Proof. Let $\mathcal{D}'_m = (Q'_m, \Sigma', \delta', 0', F')$ and $\mathcal{D}_n = (Q_n, \Sigma, \delta, 0, F)$ be minimal DFAs for L'_m and L_n , respectively. To calculate an upper bound for the boolean operations assume that $\Sigma' \setminus \Sigma$ and $\Sigma \setminus \Sigma'$ are non-empty. We add an empty state to \mathcal{D}'_m to send all transitions under the letters from $\Sigma \setminus \Sigma'$ to that state; thus we get an $(m+1)$ -state DFA $\mathcal{D}'_{m,\emptyset}$. Similarly, we add an empty state to \mathcal{D}_n to get $\mathcal{D}_{n,\emptyset}$. Now we have two DFAs over the same alphabet, and an ordinary problem of finding an upper bound for the boolean operations on two languages over the same alphabet, *except that these languages both contain empty quotients*. It is clear that $(m+1)(n+1)$ is an upper bound for all four operations; however, this bound can be improved for difference and intersection. Consider the direct product $\mathcal{P}_{m,n}$ of $\mathcal{D}'_{m,\emptyset}$ and $\mathcal{D}_{n,\emptyset}$. For difference, all $n+1$ states of $\mathcal{P}_{m,n}$ that have the form (\emptyset, q) , where $q \in Q_n$ are empty. Hence the bound can be reduced by n states to $mn + m + 1$. For intersection, all n states (\emptyset, q) , $q \in Q_n$, and all m states (p', \emptyset) , $p' \in Q'_m$, are equivalent to the empty state (\emptyset, \emptyset) , thus reducing the upper bound to $mn + 1$.

To prove that the bounds are tight, we start with $\mathcal{D}_n(a, b, c, d)$ of Definition 1. For $m, n \geq 3$, let $\mathcal{D}'_m(a, b, -, c)$ be the dialect of $\mathcal{D}'_m(a, b, c, d)$ where c plays the role of d and the alphabet is restricted to $\{a, b, c\}$, and let $\mathcal{D}_n(b, a, -, d)$ be the dialect of $\mathcal{D}_n(a, b, c, d)$ in which a and b are permuted, and the alphabet is restricted to $\{a, b, d\}$; see Figure 4.

To finish the proof, we complete the two DFAs by adding empty states, and construct their direct product as illustrated in Figure 5. If we restrict both DFAs to the alphabet $\{a, b\}$, we have the usual problem of determining the complexity of two DFAs over the same alphabet. By [1, Theorem 1], all mn states of the form $\{p', q\}$, $p' \in Q'_m$, $q \in Q_n$, are reachable and pairwise distinguishable by words in $\{a, b\}^*$ for all proper boolean operations if $(m, n) \notin \{(3, 4), (4, 3), (4, 4)\}$. For our application, the three exceptional cases were verified by computation.

To prove that the remaining states are reachable, observe that $(0', 0) \xrightarrow{d} (\emptyset', 0)$ and $(\emptyset', 0) \xrightarrow{b^q} (\emptyset', q)$, for $q \in Q_n$. Symmetrically, $(0', 0) \xrightarrow{c} (0', \emptyset)$ and $(0', \emptyset) \xrightarrow{a^{p'}} (p', \emptyset)$, for $p' \in Q'_m$. Finally, $(\emptyset', n-1) \xrightarrow{c} (\emptyset', \emptyset)$, and all $(m+1)(n+1)$ states of the direct product are reachable.

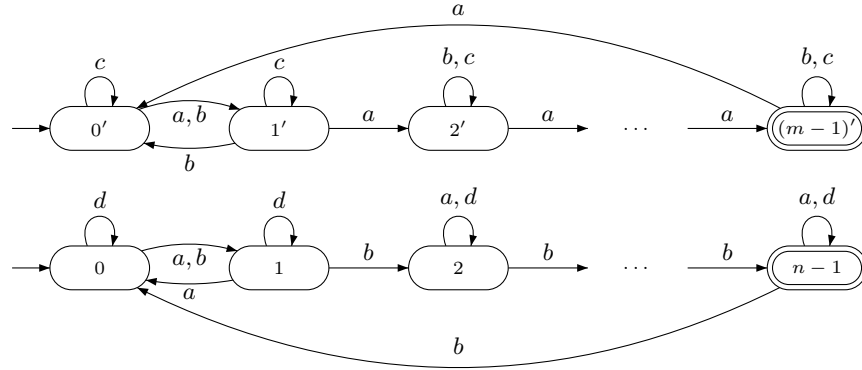


Fig. 4. Witnesses $D'_m(a, b, -, c)$ and $D_n(b, a, -, d)$ for boolean operations.

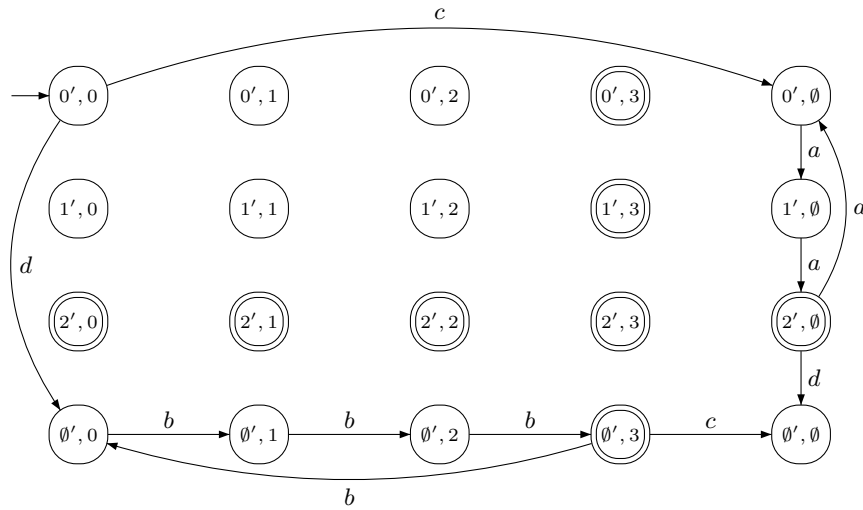


Fig. 5. Direct product for union shown partially.

It remains to verify that the appropriate states are pairwise distinguishable. From [1, Theorem 1], we know that all states in $Q'_m \times Q_n$ are distinguishable. Let $H = \{(\emptyset', q) \mid q \in Q_n\}$, and $V = \{(p', \emptyset) \mid p' \in Q'_m\}$. For the operations consider four cases:

Union The final states of $\mathcal{P}_{m,n}$ are $\{((m-1)', q) \mid q \in Q_n \cup \{\emptyset\}\}$, and $\{(p', n-1) \mid p' \in Q'_m \cup \{\emptyset'\}\}$. Every state in V accepts a word with a c , whereas no state in H accepts such words. Similarly, every state in H accepts a word with a d , whereas no state in V accepts such words. Every state in $Q'_m \times Q_n$ accepts a word with a c and a word with a d . State (\emptyset', \emptyset) accepts no words at all. Hence any two states chosen from different sets (the sets being $Q'_m \times Q_n$, H , V , and $\{(\emptyset', \emptyset)\}$) are distinguishable. States in H are distinguishable by

words in b^* and those in V , by words in a^* . Therefore all $mn + m + n + 1$ states are pairwise distinguishable.

Symmetric Difference The final states here are all the final states for union except $((m - 1)', n - 1)$. The rest of the argument is the same as for union.

Difference The final states now are $\{((m - 1)', q) \mid q \neq n - 1\}$. The n states of the form (\emptyset', q) , $q \in Q_n$, are now equivalent to the empty state $\{(\emptyset', \emptyset)\}$. The remaining states are pairwise distinguishable by the arguments used for union. Hence we have $mn + m + 1$ distinguishable states.

Intersection Here only $((m - 1)', n - 1)$ is final and all states (p', \emptyset) , $p' \in Q'_m$, and (\emptyset', q) , $q \in Q_n$ are equivalent to $\{(\emptyset', \emptyset)\}$, leaving $mn + 1$ distinguishable states. \square

Remark 1 (Marek Szykuła, personal communication). In the case of intersection the alphabet of one of the witnesses can be binary: $L'_m(a, b, -, c)$ and $L_n(b, a)$ meet the bound $mn + 1$. Reachability and distinguishability of all mn states of the form $\{p', q\}$, $p' \in Q'_m$, $q \in Q_n$, is the same as above. State (p', \emptyset) can be reached from $(p', 0)$ by c , and is equivalent to the empty state, thus giving $mn + 1$ states in the intersection.

4 Product

Theorem 2. For $m, n \geq 3$, let L'_m (respectively, L_n) be a regular language with m (respectively, n) quotients over an alphabet Σ' , (respectively, Σ). Then $\kappa(L'_m L_n) = m2^n + 2^{n-1}$.

Proof. First we derive the upper bound. Let $\mathcal{D}'_m = (Q'_m, \Sigma', \delta', 0', F')$ and $\mathcal{D}_n = (Q_n, \Sigma, \delta, 0, F)$ be minimal DFAs of L'_m and L_n , respectively. We use the normal construction of an ε -NFA \mathcal{N} to recognize $L'_m L_n$, by introducing an ε -transition from each final state of \mathcal{D}'_m to the initial state of \mathcal{D}_n , and changing all final states of \mathcal{D}'_m to non-final. This is illustrated in Figure 6, where $(m - 1)'$ is the only final state of \mathcal{D}'_m . We then determinize \mathcal{N} using the subset construction to get the DFA \mathcal{D} for $L'_m L_n$.

Suppose \mathcal{D}'_m has k final states, where $1 \leq k \leq m - 1$. I will show that \mathcal{D} can have only the following types of states: (a) at most $(m - k)2^n$ states $\{p'\} \cup S$, where $p' \in Q'_m \setminus F'$, and $S \subseteq Q_n$, (b) at most $k2^{n-1}$ states $\{p', 0\} \cup S$, where $p' \in F'$ and $S \subseteq Q_n \setminus \{0\}$, and (c) at most 2^n states $S \subseteq Q_n$. Because \mathcal{D}'_m is deterministic, there can be at most one state p' of Q'_m in any reachable subset. If $p' \notin F'$, it may be possible to reach any subset of states of Q_n along with p' , and this accounts for (a). If $p' \in F'$, then the set must contain 0 and possibly any subset of $Q_n \setminus \{0\}$, giving (b). It may also be possible to have any subset S of Q_n by applying an input that is not in Σ' to $\{0'\} \cup S$ to get S , and so we have (c). Altogether, there are at most $(m - k)2^n + k2^{n-1} + 2^n = (2m - k)2^{n-1} + 2^n$ reachable subsets. This expression reaches its maximum when $k = 1$, and hence we have at most $m2^n + 2^{n-1}$ states in \mathcal{D} .

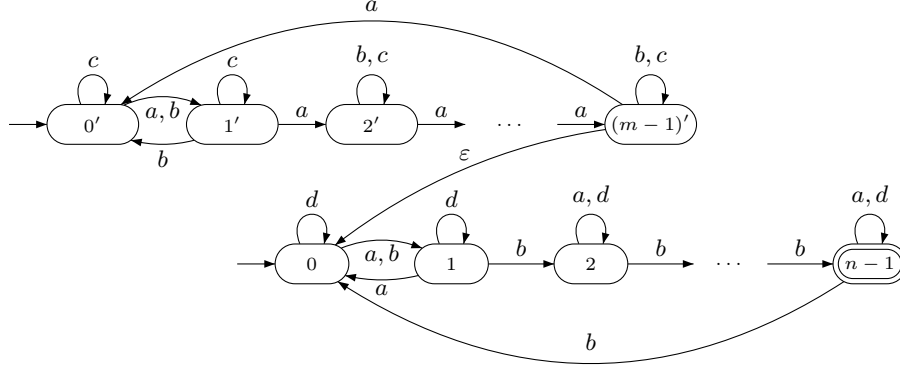


Fig. 6. An NFA for the product of $L'_m(a, b, -, c)$ and $L_n(b, a, -, d)$.

To prove that the bound is tight, we use the same witnesses as for boolean operations; see Figure 6. If $S = \{q_1, \dots, q_k\} \subseteq Q_n$ then $S+i = \{q_1+i, \dots, q_k+i\}$ and $S-i = \{q_1-i, \dots, q_k-i\}$, where addition and subtraction are modulo n . Note that b^2 and a^m (a^2 and b^n) act as the identity on Q'_m (Q_n). If $p < m-1$, then $\{p'\} \cup S \xrightarrow{b^2} \{p'\} \cup (S+2)$, for all $S \subseteq Q_n$. If n is odd, then $(b^2)^{(n-1)/2} = b^{n-1}$ and $\{p'\} \cup S \xrightarrow{b^{n-1}} \{p'\} \cup (S-1)$, for all $q \in Q_n$. If $0, 1 \notin S$ or $\{0, 1\} \subseteq S$, then a acts as the identity on S .

Remark 2. If $1 \notin S$ and $\{(m-2)'\} \cup S$ is reachable, then $\{0', 1\} \cup S$ is reachable for all $S \subseteq Q_n \setminus \{1\}$.

Proof. If $0 \in S$, then $\{(m-2)', 0\} \cup S \setminus \{0\} \xrightarrow{a} \{(m-1)', 0, 1\} \cup S \setminus \{0\} \xrightarrow{a} \{0', 0, 1\} \cup S \setminus \{0\} = \{0', 1\} \cup S$. If $0 \notin S$, then $\{(m-2)'\} \cup S \xrightarrow{a} \{(m-1)', 0\} \cup S \xrightarrow{a} \{0', 1\} \cup S$. \square

We now prove that the languages of Figure 6 meet the upper bound.

Claim 1: All sets of the form $\{p'\} \cup S$, where $p' \in Q'_{m-1}$ and $S \subseteq Q_n$, are reachable. We show this by induction on the size of S .

Basis: $|S| = 0$. The initial set is $\{0'\}$, and from $\{0'\}$ we reach $\{p'\}$, $p' \in Q'_{m-1}$, by a^p , without reaching any states of Q_n . Thus the claim holds if $|S| = 0$.

Induction Assumption: $\{p'\} \cup S$, where $p' \in Q'_{m-1}$ and $S \subseteq Q_n$, is reachable if $|S| \leq k$.

Induction Step: We prove that if $|S| = k+1$, then $\{p'\} \cup S$ is reachable. Let $S = \{q_0, q_1, \dots, q_k\}$, where $0 \leq q_0 < q_1 < \dots < q_k \leq n-1$. Suppose $q \in S$. By assumption, sets $\{p'\} \cup (S \setminus \{q\} - (q-1))$ are reachable for all $p' \in Q'_{m-1}$.

• All sets of the form $\{0'\} \cup S$ are reachable.

Note that $1 \notin (S \setminus \{q\} - (q-1))$. By assumption, $\{(m-2)'\} \cup (S \setminus \{q\} - (q-1))$ is reachable. By Remark 2, $\{0', 1\} \cup (S \setminus \{q\} - (q-1))$ is reachable.

1. If there is an odd state q in S , then $\{0', 1\} \cup (S \setminus \{q\} - (q-1)) \xrightarrow{b^{q-1}} \{0', q\} \cup (S \setminus \{q\}) = \{0'\} \cup S$.
2. If there is no odd state in S and n is odd, then $S \subseteq \{0, 2, \dots, n-1\}$. Pick $q \in S$. Then $\{0', 1\} \cup (S \setminus \{q\} - (q-1)) \xrightarrow{b^q} \{0', q+1\} \cup (S \setminus \{q\} + 1) \xrightarrow{b^{n-1}} \{0', q\} \cup S \setminus \{q\} = \{0'\} \cup S$.
3. If there is no odd state and n is even, then $S \subseteq \{0, 2, \dots, n-2\}$ (so $n-1 \notin S$).
 - (a) If $0 \notin S$, then $0, 1 \notin S+1$. By 1, $\{0'\} \cup (S+1)$ is reachable, since $S+1$ contains an odd state. Then $\{0'\} \cup (S+1) \xrightarrow{a} \{1'\} \cup (S+1) \xrightarrow{b^{n-1}} \{0'\} \cup S$.
 - (b) If $2 \notin S$, then $0, 1 \notin S-1$. By 1, $\{0'\} \cup (S-1)$ is reachable, since $S-1$ contains an odd state. Then $\{0'\} \cup (S-1) \xrightarrow{a} \{1'\} \cup (S-1) \xrightarrow{b^{n+1}} \{0'\} \cup S$.
 - (c) If $\{0, 2\} \subseteq S$, then $0 \notin S-1$, and $1, n-1 \in S-1$. By 1, $\{0'\} \cup (S-1)$ is reachable, since $1 \in S-1$. Note that aba sends 1 to 0, $n-1$ to 1, and adds 1 to each state $q \geq 3$ of $S-1$; thus $2 \notin (S-1)aba$, and $\{0'\} \cup (S-1) \xrightarrow{aba} \{1', 0, 1\} \cup S \setminus \{0, 2\}$. Next, b^{n-1} sends 0 to $n-1$ and subtracts 1 from every other element of $S \setminus \{0, 2\}$. Hence $\{1', 0, 1\} \cup S \setminus \{0, 2\} \xrightarrow{b^{n-1}} \{0', n-1, 0\} \cup (S \setminus \{0, 2\} - 1) \xrightarrow{ab} \{0', 0, 2\} \cup (S \setminus \{0, 2\}) = \{0'\} \cup S$.

• All sets of the form $\{1'\} \cup S$ are reachable.

If 0 and 1 are not in S or are both in S , then $\{0'\} \cup S \xrightarrow{a} \{1'\} \cup S$. If $0 \in S$ but $1 \notin S$, then $\{0', 1\} \cup S \setminus \{0\} \xrightarrow{a} \{1', 0\} \cup S \setminus \{0\} = \{1'\} \cup S$. If $1 \in S$ but $0 \notin S$, then $\{0', 0\} \cup S \setminus \{1\} \xrightarrow{a} \{1', 1\} \cup S \setminus \{1\} = \{1'\} \cup S$.

• All sets of the form $\{p'\} \cup S$, where $2 \leq p \leq m-2$, are reachable.

If p is even, then $\{0'\} \cup S \xrightarrow{a^p} \{p'\} \cup S$.

If p is odd, then $\{1'\} \cup S \xrightarrow{a^{p-1}} \{p'\} \cup S$.

Claim 2: All sets of the form $\{(m-1)', 0\} \cup S$ are reachable.

1. By Claim 1, $\{(m-3)'\} \cup S$ is reachable. If $q_0 = 1$, then $\{(m-3)', 1\} \cup S \setminus \{1\} \xrightarrow{a^2} \{(m-1)', 0, 1\} \cup S \setminus \{1\} = \{(m-1)', 0\} \cup S$.
2. By Claim 1, $\{(m-2)'\} \cup S$ is reachable. If $q_0 \geq 2$, then $\{(m-2)'\} \cup S \xrightarrow{a} \{(m-1)', 0\} \cup S$.

Claim 3: All sets of the form S are reachable.

By Claim 1, $\{0'\} \cup S$ is reachable for every S , and $\{0'\} \cup S \xrightarrow{d} S$.

For distinguishability, note that only state q accepts $w_q = b^{n-1-q}$ in \mathcal{D}_n . Hence, if two states of the product have different sets S and S' and $q \in S \oplus S'$, then they can be distinguished by w_q . State $\{p'\} \cup S$ is distinguished from S by $ca^{m-1-p}b^{n-1}$. If $p < q$, states $\{p'\} \cup S$ and $\{q'\} \cup S$ are distinguished as follows. Use ca^{m-1-q} to reach $\{(p+m-1-q)'\}$ from p' and $\{(m-1)'\} \cup \{0\}$ from q' . The reached states are distinguishable since they differ in their subsets of Q_n . \square

5 Most Complex Regular Languages

A *most complex* regular language stream is one that, together with some dialects, meets the complexity bounds for all boolean operations, product, star, and reversal, and has the largest syntactic semigroup and most complex atoms [3]. A most complex stream should have the smallest possible alphabet sufficient to meet all the bounds. Most complex streams are useful in systems dealing with regular languages and finite automata. One would like to know the maximal sizes of automata that can be handled by the system. In view of the existence of most complex streams, one stream can be used to test all the operations. Here we present a stream similar to that of [3] but with one added input letter that induces the identity transformation, as shown in Figure 3.

Theorem 3 (Most Complex Regular Languages). *For each $n \geq 3$, the DFA of Definition 1 is minimal and its language $L_n(a, b, c, d)$ has complexity n . The stream $(L_m(a, b, c, d) \mid m \geq 3)$ with dialect streams $(L_n(a, b, -, c) \mid n \geq 3)$ and $(L_n(b, a, -, d) \mid n \geq 3)$ is most complex in the class of regular languages. In particular, it meets all the complexity bounds below, which are maximal for regular languages. In several cases the bounds can be met with a restricted alphabet.*

1. The syntactic semigroup of $L_n(a, b, c)$ has cardinality n^n .
2. Each quotient of $L_n(a)$ has complexity n .
3. The reverse of $L_n(a, b, c)$ has complexity 2^n , and $L_n(a, b, c)$ has 2^n atoms¹.
4. For each atom A_S of $L_n(a, b, c)$, the complexity $\kappa(A_S)$ satisfies:

$$\kappa(A_S) = \begin{cases} 2^n - 1, & \text{if } S \in \{\emptyset, Q_n\}; \\ 1 + \sum_{x=1}^{|S|} \sum_{y=1}^{n-|S|} \binom{n}{x} \binom{n-x}{y}, & \text{if } \emptyset \subsetneq S \subsetneq Q_n. \end{cases}$$

5. The star of $L_n(a, b)$ has complexity $2^{n-1} + 2^{n-2}$.
6. The product $L'_m(a, b, -, c)L_n(b, a, -, d)$ has complexity $m2^n + 2^{n-1}$.
7. The complexity of $L'_m(a, b, -, c) \circ L_n(b, a, -, d)$ is $mn + m + n + 1$ if $\circ \in \{\cup, \oplus\}$, $mn + m + 1$ if $\circ = \setminus$, and $mn + 1$ if $\circ = \cap$.

Proof. The proofs of 1–5 can be found in [3], and Claims 6 and 7 are proved in the present paper, Theorems 1 and 2. \square

Proposition 1 (Marek Szykuła, personal communication). *At least four inputs are required for a most complex regular language. In particular, four inputs are needed for union: two inputs are needed to reach all pairs of states in $Q'_m \times Q_n$, one input in $\Sigma' \setminus \Sigma$ for pairs (p', \emptyset) with $p' \in Q'_m$, and one in $\Sigma \setminus \Sigma'$ for pairs (\emptyset', q) with $q \in Q_n$.*

¹ The *atom congruence* is a left congruence defined as follows: two words x and y are equivalent if $ux \in L$ if and only if $uy \in L$ for all $u \in \Sigma^*$. Thus x and y are equivalent if $x \in u^{-1}L$ if and only if $y \in u^{-1}L$. An equivalence class of this relation is called an *atom* of L [5, 7]. It follows that an atom is a non-empty intersection of complemented and uncomplemented quotients of L . The number of atoms and their quotient complexities are possible measures of complexity of regular languages [3]. For more information about atoms and their complexity, see [4, 5, 7].

6 Conclusions

Two complete DFAs over different alphabets Σ' and Σ are incomplete DFAs over $\Sigma' \cup \Sigma$. Each DFA can be completed by adding an empty state and sending all transitions induced by letters not in the DFA's alphabet to that state. This results in an $(m + 1)$ -state DFA and an $(n + 1)$ -state DFA. From the theory about DFAs over the same alphabet we know that $(m + 1)(n + 1)$ is an upper bound for all boolean operations on the original DFAs, and that $m2^{n+1} + 2^n$ is an upper bound for product. We have shown that the tight bounds for boolean operations are $(m + 1)(n + 1)$ for union and symmetric difference, $mn + m + 1$ for difference, and $mn + 1$ for intersection, while the tight bound for product is $m2^n + 2^{n-1}$. In the same-alphabet case the tight bound is mn for all boolean operations and it is $(m - 1)2^n + 2^{n-1}$ for product. In summary, the restriction of identical alphabets is unnecessary and leads to incorrect results.

It should be noted that if the two languages in question already have empty quotients, then making the alphabets the same does not require the addition of any states, and the traditional same-alphabet methods are correct. This is the case, for example, for prefix-free, suffix-free and finite languages.

Acknowledgment I am very grateful to Sylvie Davies, Bo Yang Victor Liu and Corwin Sinnamon for careful proofreading and constructive comments. I thank Marek Szykula for contributing the important Remark 1 and Proposition 1.

References

1. Bell, J., Brzozowski, J., Moreira, N., Reis, R.: Symmetric groups and quotient complexity of boolean operations. In: Esparza, J., et al. (eds.) ICALP 2014. LNCS, vol. 8573, pp. 1–12. Springer (2014)
2. Brzozowski, J.: Quotient complexity of regular languages. *J. Autom. Lang. Comb.* 15(1/2), 71–89 (2010)
3. Brzozowski, J.: In search of the most complex regular languages. *Int. J. Found. Comput. Sc.* 24(6), 691–708 (2013)
4. Brzozowski, J., Tamm, H.: Complexity of atoms of regular languages. *Int. J. Found. Comput. Sc.* 24(7), 1009–1027 (2013)
5. Brzozowski, J., Tamm, H.: Theory of atomata. *Theoret. Comput. Sci.* 539, 13–27 (2014)
6. Gao, Y., Salomaa, K., Yu, S.: Transition complexity of incomplete DFAs. *Fund. Inform.* 110, 143–158 (2011)
7. Iván, S.: Complexity of atoms, combinatorially. *Inform. Process. Lett.* 116(5), 356–360 (2016)
8. Maia, E., Moreira, N., Reis, R.: Incomplete operational transition complexity of regular languages. *Inform. and Comput.* 244, 1–22 (2015)
9. Maslov, A.N.: Estimates of the number of states of finite automata. *Dokl. Akad. Nauk SSSR* 194, 1266–1268 (Russian). (1970), English translation: *Soviet Math. Dokl.* 11 (1970) 1373–1375
10. Yu, S.: State complexity of regular languages. *J. Autom. Lang. Comb.* 6, 221–234 (2001)
11. Yu, S., Zhuang, Q., Salomaa, K.: The state complexities of some basic operations on regular languages. *Theoret. Comput. Sci.* 125, 315–328 (1994)