



STATE COMPLEXITY OF MULTIPLE CONCATENATION

JOZEF JIRÁSEK ^(A,B) GALINA JIRÁSKOVÁ ^(C,D)

^(A) *Institute of Computer Science, P. J. Šafárik University*
Jesenná 5, 040 01 Košice, Slovakia
 jozef.jirasek@upjs.sk

^(C) *Mathematical Institute, Slovak Academy of Sciences*
Grešákova 6, 040 01 Košice, Slovakia
 jiraskov@saske.sk

ABSTRACT

We describe witness languages meeting the upper bound on the state complexity of the multiple concatenation of k regular languages over an alphabet of size $k + 1$ with a significantly simpler proof than that in the literature. We also consider the case where some languages may be recognized by two-state automata. Then we show that one symbol can be saved, and we define witnesses for the multiple concatenation of k languages over a k -letter alphabet. This solves an open problem stated by Caron et al. [2018, Fundam. Inform. 160, 255–279]. We prove that for the concatenation of three languages, the ternary alphabet is optimal. We also show that a trivial upper bound on the state complexity of multiple concatenation is asymptotically tight for ternary languages, and that a lower bound remains exponential in the binary case. Finally, we obtain a tight upper bound for unary cyclic languages and languages recognized by unary automata that do not have final states in their tails.

Keywords: regular languages, multiple concatenation, state complexity


1. Introduction

Given formal languages L_1, L_2, \dots, L_k over an alphabet Σ , their concatenation is the language $L_1 L_2 \cdots L_k = \{u_1 u_2 \cdots u_k \mid u_i \in L_i \text{ for } i = 1, 2, \dots, k\}$. Here we consider the case where all languages are regular and ask the question of how many states are sufficient and necessary in the worst case for a deterministic finite automaton to recognize their concatenation assuming that each L_i is recognized by an n_i -state deterministic finite automaton.

A preliminary version of this paper appeared in Proc. DCFS 2020, LNCS vol. 12442, pp. 78–90.

^(B) Research supported by VEGA grant 1/0350/22.

^(D) Research supported by VEGA grant 2/0096/23.

 Jozef Jirásek: 0000-0003-4822-230X, Galina Jirásková: 0000-0001-9817-8197

The first results for the concatenation of two regular languages were obtained by Maslov [5] in 1970. In particular, he described binary witnesses meeting the upper bound $n_1 2^{n_2} - 2^{n_2-1}$. In 1994 Yu et al. [8] proved that this upper bound cannot be met if the first language is recognized by a minimal deterministic finite automaton that has more than one final state.

The concatenation of three and four regular languages was considered by Ésik et al. [2] in 2009, where the witnesses for the concatenation of three languages over a five-letter alphabet can be found. The rather complicated expression for the upper bounds for the concatenation of k languages, as well as witnesses over a $(2k-1)$ -letter alphabet were given by Gao and Yu [4].

Caron et al. [1] presented recursive formulas for the upper bounds, and described witnesses over a $(k+1)$ -letter alphabet using Brzozowski's universal automata. They also showed that to meet the upper bound for the concatenation of two or three languages, the binary or ternary alphabet, respectively, is enough, and they conjectured that k symbols could be enough to describe witnesses for the concatenation of k languages.

In this paper, we study in detail the state complexity of multiple concatenation of k regular languages. We first describe witnesses over an alphabet consisting of $k+1$ symbols with a significantly simpler proof than that in [1]. Our witness automata A_1, A_2, \dots, A_k are defined over the alphabet $\{b, a_1, \dots, a_k\}$. Each a_i performs the circular shift in A_i and the identity in all the other automata. These k permutation symbols are used to get the reachability of all so-called valid states in a DFA for concatenation. The symbol b performs a contraction in each A_i and assures the distinguishability of all valid states almost for free. However, the proof requires that each A_i has at least three states. With a slightly more complicated proof, we also solve the case that includes two-state automata. Then we describe special binary witnesses for the concatenation of two languages. We combine our ideas used for the $(k+1)$ -letter alphabet and those for binary witnesses to describe witnesses for multiple concatenation over a k -letter alphabet, which solves an open problem stated by Caron et al. [1]. In the case of $k=3$, we show that the ternary alphabet is optimal.

We also examine multiple concatenation on binary, ternary, and unary languages. We show that in the binary case, the lower bounds remain exponential in n_2, n_3, \dots, n_k , and in the ternary case, the trivial upper bound $n_1 2^{n_2+n_3+\dots+n_k}$ can be met up to some multiplicative constant depending on k . For unary languages, we use Frobenius numbers to get a tight upper bound for cyclic languages, or languages recognized by automata that do not have final states in their tails. We also consider the case with final states in tails, and provide upper and lower bounds for multiple concatenation in such a case.

2. Preliminaries

We assume that the reader is familiar with basic notions in automata and formal language theory. For details and all unexplained notions, we refer the reader to [7]. The size of a finite set S is denoted by $|S|$, and the set of all its subsets by 2^S .

For a finite non-empty alphabet of symbols Σ , the set of all strings over Σ , including the empty string ε , is denoted by Σ^* . A language is any subset of Σ^* . The multiple concatenation of k languages L_1, L_2, \dots, L_k is the language $L_1 L_2 \cdots L_k = \{u_1 u_2 \cdots u_k \mid u_1 \in L_1, u_2 \in L_2, \dots, u_k \in L_k\}$.

A *deterministic finite automaton* (DFA) is a quintuple $A = (Q, \Sigma, \cdot, s, F)$ where Q is a non-empty finite set of states, Σ is a non-empty finite alphabet of input symbols, $\cdot: Q \times \Sigma \rightarrow Q$ is the transition function, $s \in Q$ is the initial state, and $F \subseteq Q$ is the set of final (accepting) states. The transition function can be naturally extended to the domain $Q \times \Sigma^*$. The language recognized (accepted) by the DFA A is the set of strings $L(A) = \{w \in \Sigma^* \mid s \cdot w \in F\}$.

All deterministic finite automata in this paper are assumed to be complete; that is, the transition function is a total function.

We usually omit \cdot , and write qa instead of $q \cdot a$. Next, for a subset S of Q and a string w , let $Sw = \{qw \mid q \in S\}$ and $wS = \{q \mid qw \in S\}$. Each input symbol a induces a transformation on $Q = \{q_1, q_2, \dots, q_n\}$ given by $q \mapsto qa$. We denote by $a: (q_1, q_2, \dots, q_\ell)$ the transformation that maps q_i to q_{i+1} for $i = 1, \dots, \ell - 1$, the state q_ℓ to q_1 , and fixes any other state in Q . In particular, (q_1) denotes the identity. Next, we denote by $a: (q_1 \rightarrow q_2 \rightarrow \cdots \rightarrow q_\ell)$ the transformation that maps q_i to q_{i+1} for $i = 1, 2, \dots, \ell - 1$ and fixes any other state. Finally, we denote by $a: (S \rightarrow q_i)$ the transformation that maps each $q \in S$ to q_i and fixes any other state.

A state $q \in Q$ is *reachable* in the DFA A if there is a string $w \in \Sigma^*$ such that $q = sw$. Two states p and q are *distinguishable* if there is a string w such that exactly one of the states pw and qw is final. A state $q \in Q$ is a *dead state* if $qw \notin F$ for every string $w \in \Sigma^*$.

A DFA is *minimal* (with respect to the number of states) if all its states are reachable and pairwise distinguishable. The *state complexity* of a regular language L , $\text{sc}(L)$, is the number of states in the minimal DFA recognizing L . The state complexity of a k -ary regular operation f is a function from \mathbb{N}^k to \mathbb{N} given by $(n_1, n_2, \dots, n_k) \mapsto \max\{\text{sc}(f(L_1, L_2, \dots, L_k)) \mid \text{sc}(L_i) \leq n_i \text{ for } i = 1, 2, \dots, k\}$.

A *nondeterministic finite automaton* (NFA) is a quintuple $N = (Q, \Sigma, \cdot, I, F)$ where Q, Σ , and F are the same as for a DFA, $I \subseteq Q$ is the set of initial states, and $\cdot: Q \times (\Sigma \cup \{\varepsilon\}) \rightarrow 2^Q$ is the transition function. A string w in Σ^* is *accepted* by the NFA N if $w = a_1 a_2 \cdots a_m$ where $a_i \in \Sigma \cup \{\varepsilon\}$ and a sequence of states q_0, q_1, \dots, q_m exists in Q such that $q_0 \in I$, $q_{i+1} \in q_i \cdot a_{i+1}$ for $i = 0, 1, \dots, m - 1$, and $q_m \in F$. The language recognized by the NFA N is the set of strings $L(N) = \{w \in \Sigma^* \mid w \text{ is accepted by } N\}$. For $p, q \in Q$ and $a \in \Sigma \cup \{\varepsilon\}$, we say that a triple (p, a, q) is a *transition* in N if $q \in p \cdot a$.

Let $N = (Q, \Sigma, \cdot, I, F)$ be an NFA. For a set $S \subseteq Q$, let $E(S)$ denote the ε -closure of S ; that is, the set of states $\{q \mid q \text{ is reached from a state in } S \text{ through 0 or more } \varepsilon\text{-transitions}\}$. The *subset automaton* of the NFA N is the DFA $\mathcal{D}(N) = (2^Q, \Sigma, \cdot', E(I), F')$ where $F' = \{S \in 2^Q \mid S \cap F \neq \emptyset\}$ and $S \cdot' a = \cup_{q \in S} E(q \cdot a)$ for each $S \in 2^Q$ and each $a \in \Sigma$. The subset automaton $\mathcal{D}(N)$ recognizes the language $L(N)$.

The *reverse* of the NFA N is the NFA $N^R = (Q, \Sigma, \cdot^R, F, I)$ where the transition function is defined by $q \cdot^R a = \{p \in Q \mid q \in p \cdot a\}$; that is, N^R is obtained from N by

swapping the roles of initial and final states, and by reversing all transitions.

A subset S of Q is *reachable* in N if there is a string w in Σ^* such that $S = I \cdot w$, and it is *co-reachable* in N if it is reachable in the reverse N^R .

We use the following two simple observations to prove distinguishability of states in subset automata.

Lemma 1. *Let $N = (Q, \Sigma, \cdot, I, F)$ be an NFA without ε -transitions. Let $S, T \subseteq Q$ and $q \in S \setminus T$. If the singleton set $\{q\}$ is co-reachable in N , then S and T are distinguishable in the subset automaton $\mathcal{D}(N)$.*

Proof. Since the singleton set $\{q\}$ is co-reachable in N , there is a string $w \in \Sigma^*$ which sends the set of final states F to $\{q\}$ in the reversed automaton N^R . It follows that the string w^R is accepted by N from the state q , and it is rejected from any other state. Thus, the string w^R is accepted by $\mathcal{D}(N)$ from S and rejected from T . \square

Corollary 2. *If for each state q of an NFA N , the singleton set $\{q\}$ is co-reachable in N , then all states of the subset automaton $\mathcal{D}(N)$ are pairwise distinguishable.* \square

3. Multiple Concatenation: Upper Bound

In this section, we recall the constructions of ε -NFAs and NFAs for multiple concatenation, as well as the known upper bounds. We also provide a simple alternative method to get upper bounds. In the last part of this section, we consider the case when some of given automata have just one state.

For $i = 1, 2, \dots, k$, let $A_i = (Q_i, \Sigma, \cdot_i, s_i, F_i)$ be a DFA, and assume that $Q_i \cap Q_j = \emptyset$ if $i \neq j$. Then the concatenation $L(A_1)L(A_2) \cdots L(A_k)$ is recognized by an NFA $N = (Q_1 \cup Q_2 \cup \dots \cup Q_k, \Sigma, \cdot, s_1, F_k)$, where for each $i = 1, 2, \dots, k$, each $q \in Q_i$, and each $a \in \Sigma$, we have $q \cdot a = \{q \cdot_i a\}$ and for each $i = 1, 2, \dots, k-1$ and each $q \in F_i$, we have $q \cdot \varepsilon = \{s_{i+1}\}$, that is, the NFA N is obtained from the DFAs A_1, A_2, \dots, A_k by adding the ε -transition from each final state of A_i to the initial state s_{i+1} of A_{i+1} for $i = 1, 2, \dots, k-1$; the initial state of N is s_1 , and its set of final states is F_k .

Since A_1 is a complete DFA, in the corresponding subset automaton $\mathcal{D}(N)$, each reachable subset is of the form $\{q\} \cup S_2 \cup S_3 \cup \dots \cup S_k$ where $q \in S_1$ and $S_i \subseteq Q_i$ for $i = 2, 3, \dots, k$. We represent such a set by the k -tuple $(\{q\}, S_2, S_3, \dots, S_k)$, or more often by $(q, S_2, S_3, \dots, S_k)$, and with this representation, it is not necessary to have the state sets disjoint. Nevertheless, since we sometimes use special properties of the NFA N , we keep in mind that this k -tuple represents the union of appropriate set of states of the corresponding DFAs. We usually denote all transition functions by \cdot , and simply write $(qa, S_2, S_3, \dots, S_k)$ or $(q, S_2a, S_3, \dots, S_k)$; that is, applying a to the i -th component means that we use the transition function \cdot_i .

It follows from the construction of the NFA N that if $S_i \cap F_i \neq \emptyset$ then $s_{i+1} \in S_{i+1}$, and if $S_i = \emptyset$, then $S_{i+1} = \emptyset$ in any reachable state (S_1, S_2, \dots, S_k) of the subset automaton $\mathcal{D}(N)$. The states satisfying the above mentioned properties are called valid in [1]; let us summarize the three properties in the next definition.

Definition 3. A state (S_1, S_2, \dots, S_k) of the subset automaton $\mathcal{D}(N)$ is *valid* if

- (I) $|S_1| = 1$,
- (II) if $S_i = \emptyset$ and $i \leq k-1$, then $S_{i+1} = \emptyset$,
- (III) if $S_i \cap F_i \neq \emptyset$ and $i \leq k-1$, then $s_{i+1} \in S_{i+1}$.

Since each reachable state of $\mathcal{D}(N)$ is valid, we have the next observation.

Proposition 4. An upper bound on $\text{sc}(L(A_1)L(A_2) \cdots L(A_k))$ is given by the number of valid states in the subset automaton $\mathcal{D}(N)$. \square

Notice that, to reach as many valid states as possible, each automaton A_i with $i \leq k-1$ should have exactly one final state f_i , that is, we have $F_i = \{f_i\}$. Moreover, if A_i has at least two states, then we should have $s_i \neq f_i$. If this is the case for all A_i , then we can construct an NFA N for the concatenation $L(A_1)L(A_2) \cdots L(A_k)$ from the DFAs A_1, A_2, \dots, A_k as follows: for each $i = 1, 2, \dots, k-1$, each state $q \in Q_i$, and each symbol $a \in \Sigma$ such that $q \cdot_i a = f_i$, we add the transition (q, a, s_{i+1}) ; the initial state of N is s_1 , and its unique final state is f_k .

For $k = 2$, an upper bound on the number of valid states is $(n_1 - 1)2^{n_2} + 2^{n_2-1}$ [8], which is the sum of the number of states (q, S_2) with $q \neq f_1$ and $S_2 \subseteq Q_2$ and the number of states (f_1, S_2) with $s_2 \in S_2$. For $k \geq 3$, we have the following inequalities.

Proposition 5. Let $k \geq 3$ and $\# \tau_k$ denote the number of valid states. Then

$$\frac{1}{2^{k-1}} n_1 2^{n_2+n_3+\dots+n_k} \leq \# \tau_k \leq \frac{3}{4} n_1 2^{n_2+n_3+\dots+n_k}.$$

Proof. Every state (S_1, S_2, \dots, S_k) with $s_i \in S_i$ for $i = 2, 3, \dots, k$ is a valid state. This gives the left inequality. On the other hand, every state (S_1, S_2, \dots, S_k) with $f_2 \in S_2$ and $s_3 \notin S_3$ is not valid, which gives the right inequality. \square

We now provide a simple alternative method for obtaining an upper bound on the number of valid states. To this aim let

- U_i be the number of tuples $(S_i, S_{i+1}, \dots, S_k)$ such that for fixed $S'_1, S'_2, \dots, S'_{i-1}$ with $f_{i-1} \notin S'_{i-1}$ the state $(S'_1, \dots, S'_{i-1}, S_i, S_{i+1}, \dots, S_k)$ is valid,
- V_i be the number of tuples $(S_i, S_{i+1}, \dots, S_k)$ such that for a fixed $S'_1, S'_2, \dots, S'_{i-1}$ with $f_{i-1} \in S'_{i-1}$ the state $(S'_1, \dots, S'_{i-1}, S_i, S_{i+1}, \dots, S_k)$ is valid.

Then we have the next result.

Theorem 6. Let $k \geq 2$, $n_i \geq 2$ for $i = 1, 2, \dots, k$, and $A_i = (Q_i, \Sigma, \cdot, s_i, \{f_i\})$ be an n_i -state DFA with $s_i \neq f_i$. Let U_i and V_i be as defined above, and $\# \tau_k$ be the number of valid states in the subset automaton $\mathcal{D}(N)$ accepting $L(A_1)L(A_2) \cdots L(A_k)$. Then

$$U_k = 2^{n_k} \text{ and } V_k = 2^{n_k-1}, \quad (1)$$

and for $i = 2, 3, \dots, k-1$,

$$U_i = 1 + (2^{n_i-1} - 1)U_{i+1} + 2^{n_i-1}V_{i+1}, \quad (2)$$

$$V_i = 2^{n_i-2}(U_{i+1} + V_{i+1}). \quad (3)$$

Finally, we have

$$\# \tau_k = (n_1 - 1)U_2 + V_2. \quad (4)$$

Proof. If $f_{k-1} \notin S'_{k-1}$, then S_k may be an arbitrary subset of Q_k . If $f_{k-1} \in S'_{k-1}$, then S_k must contain s_k . This gives (1).

Let $f_{i-1} \notin S'_{i-1}$. Then we have just one tuple with $S_i = \emptyset$, namely, $(\emptyset, \emptyset, \dots, \emptyset)$, then $(2^{n_i} - 1)U_{i+1}$ tuples with $f_i \notin S_i$ and S_i non-empty, and $2^{n_i-1}V_{i+1}$ tuples with $f_i \in S_i$ final. This gives (2).

Let $f_{i-1} \in S'_{i-1}$. Then $s_i \in S_i$. We have $(2^{n_i} - 2)U_{i+1}$ tuples with $s_i \in S_i$ and $f_i \notin S_i$, and $2^{n_i-2}V_{i+1}$ tuples with $s_i \in S_i$ and $f_i \in S_i$. This gives (3).

Finally, we have $(n_1 - 1)$ possibilities for S'_1 to be non-final singleton set, and one, namely, $S'_1 = \{f_1\}$, to be final. This gives (4). \square

Let us illustrate the above result in the following example.

Example 7. Let $k = 3$ and $n_1, n_2, n_3 \geq 2$. Then

$$\begin{aligned} U_3 &= 2^{n_3} \text{ and } V_3 = 2^{n_3-1}, \\ U_2 &= 1 + (2^{n_2-1} - 1)U_3 + 2^{n_2-1}V_3 = 1 + (2^{n_2-1} - 1)2^{n_3} + 2^{n_2-1}2^{n_3-1}, \\ V_2 &= 2^{n_2-2}(U_3 + V_3) = 2^{n_2-2}(2^{n_3} + 2^{n_3-1}) \\ \# \tau_k &= (n_1 - 1)U_2 + V_2 = \\ &= (n_1 - 1)(1 + (2^{n_2-1} - 1)2^{n_3} + 2^{n_2-1}2^{n_3-1}) + 2^{n_2-2}(2^{n_3} + 2^{n_3-1}) = \\ &= n_1(1 + 2^{n_2+n_3-1} - 2^{n_3} + 2^{n_2+n_3-2}) - 1 - 2^{n_2+n_3-1} + 2^{n_3} - 2^{n_2+n_3-2} + \\ &= 2^{n_2+n_3-2} + 2^{n_2+n_3-3} = \\ &= n_1(1 + \frac{3}{4}2^{n_2+n_3} - 2^{n_3}) - \frac{3}{8}2^{n_2+n_3} + 2^{n_3} - 1, \end{aligned}$$

which is the same as in [1, Example 3.6]. \blacksquare

To conclude this section, let us consider also the case when some automata have just one state. If this state is non-final, then the resulting concatenation is empty. Thus, assume that all one-state automata recognize Σ^* , so consist of one initial and final state f_i . Then we construct an NFA N accepting the language $L(A_1)L(A_2) \cdots L(A_k)$ as described above. Let $\mathcal{D}(N)$ be the corresponding subset automaton. We represent its states by k -tuples $(\{q\}, S_2, S_3, \dots, S_k)$ where $q \in Q_1$ and $S_i \subseteq Q_i$. Moreover, if $n_i = 1$, then $S_i = \{f_i\}$. If $n_i \geq 2$ and $i < k$, then to get maximum number of valid reachable sets, we must have $F_i = \{f_i\}$ and $s_i \neq f_i$. The next observation provides an upper bound in the case when exactly one of given DFAs has one state.

Proposition 8. Let $k \geq 2$, $j \in \{1, 2, \dots, k\}$, $n_j = 1$, and $n_i \geq 2$ if $i \neq j$. For $i = 1, 2, \dots, k$, let A_i be an n_i -state DFA and $L = L(A_1)L(A_2) \cdots L(A_k)$. Let U_i and V_i be given by expressions (2)-(3). Then

$$\text{sc}(L) \leq \begin{cases} V_2, & \text{if } j = 1; \\ n_1, & \text{if } j = k = 2; \\ (n_1 - 1)U_2 + V_2 + 1 \\ \quad \text{with } U_{k-1} = 2^{n_{k-1}-1} \text{ and } V_{k-1} = 2^{n_{k-1}-2}, & \text{if } j = k \geq 3; \\ (n_1 - 1)U_2 + V_2 + V_{i+1} \\ \quad \text{with } U_{j-1} = 2^{n_{j-1}-1} \text{ and } V_{j-1} = 2^{n_{i-1}-2}, & \text{if } 2 \leq j \leq k - 1. \end{cases}$$

Proof. First, let $j = 1$. Then we have $S_1 = \{f_1\}$ in each valid state (S_1, S_2, \dots, S_k) . It follows that the number of valid states is V_2 with $U_k = 2^{n_k}$ and $V_k = 2^{n_k-1}$.

Now, let $j = k$. Then all states $(S_1, S_2, \dots, S_{k-1}, \{f_k\})$ are equivalent to a final sink state. If $S_k = \emptyset$, then $f_{k-1} \notin S_{k-1}$. This results in an upper bound $(n_1 - 1)U_2 + V_2 + 1$ with $U_{k-1} = 2^{n_{k-1}-1}$ and $V_{k-1} = 2^{n_{k-1}-2}$ if $k \geq 3$ and $(n_1 - 1) + 1$ if $k = 2$.

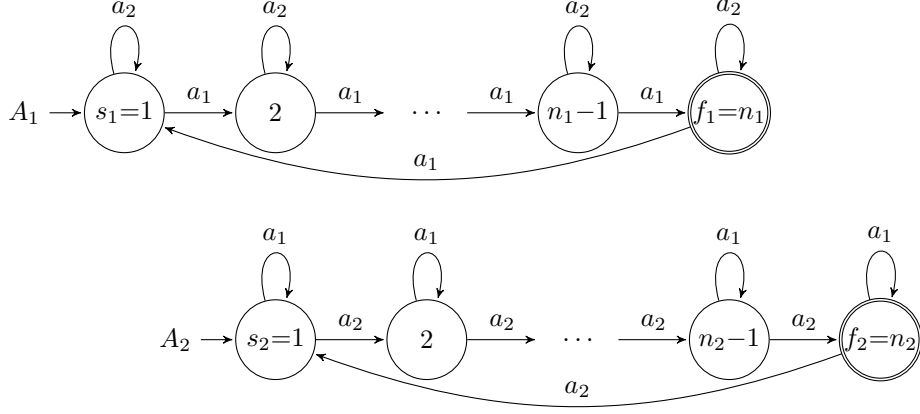
Finally, let $2 \leq j \leq k - 1$. Then all states $(S_1, S_2, \dots, S_{i-1}, \{f_i\}, \{s_{i+1}\}, \emptyset, \emptyset, \dots, \emptyset)$ are equivalent to the state $(\{s_1\}, \{s_2\}, \dots, \{s_{i-1}\}, \{f_i\}, \{s_{i+1}\}, \emptyset, \emptyset, \dots, \emptyset)$ since we have a loop on each input symbol in the state f_i and therefore every string accepted by N from a state in $Q_1 \cup Q_2 \cup \dots \cup Q_{i-1}$ is accepted also from f_i . It follows that the reachable and pairwise distinguishable valid states of $\mathcal{D}(N)$ are either of the form $(S_1, S_2, \dots, S_{i-1}, \emptyset, \emptyset, \dots, \emptyset)$ or of the form $(\{s_1\}, \{s_2\}, \dots, \{s_{i-1}\}, \{f_i\}, S_{i+1}, S_{i+2}, \dots, S_k)$. If $S_i = \emptyset$, then S_{i-1} does not contain f_i , so the number of valid states of the first form is given by $(n_i - 1)U_2 + V_2$ with $U_{i-1} = 2^{n_{i-1}-1}$ and $V_{i-1} = 2^{n_{i-1}-2}$. The number of valid states of the second form is given by V_{i+1} . \square

Example 9. Let $k = 4$, $n_3 = 1$, and $n_1, n_2, n_4 \geq 2$. Then number of valid states $(S_1, S_2, \emptyset, \emptyset)$ is $(n_1 - 1)U_2 + V_2$ where $U_2 = 2^{n_2-1}$ and $V_2 = 2^{n_2-2}$. Next, the number of valid states $(\{s_1\}, \{s_2\}, \{f_3\}, S_4)$ is $V_4 = 2^{n_4-1}$. This gives an upper bound $(n_1 - 1)2^{n_2-1} + 2^{n_2-2} + 2^{n_4-1}$ for concatenation of four languages, the third of which is Σ^* . \blacksquare

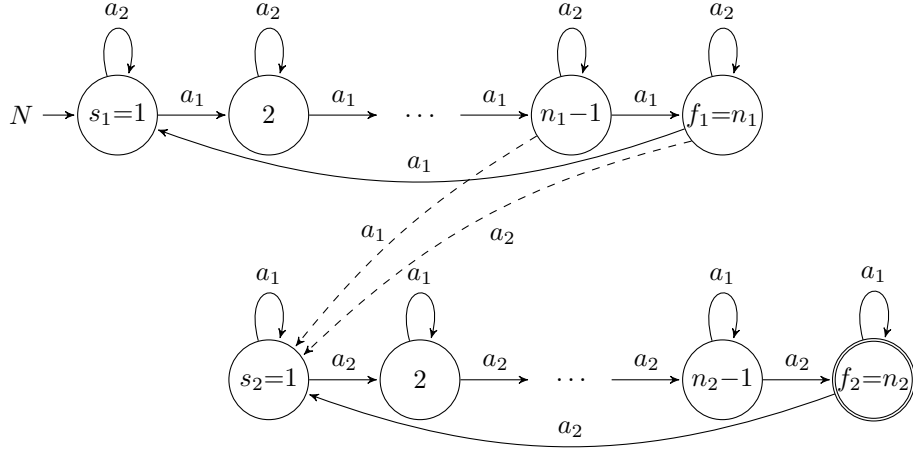
4. Matching Lower Bound: $(k + 1)$ -letter Alphabet

In this section, we describe witness languages meeting the upper bound on the state complexity of multiple concatenation of k regular languages over a $(k + 1)$ -letter alphabet with a significantly simpler proof than that in [1, Section 4, pp. 266–271]. We use these witnesses in the next section to describe witness languages over a k -letter alphabet. Let us start with the following example.

Example 10. Let $n_1, n_2 \geq 3$. Consider DFAs A_1 and A_2 over $\{a_1, a_2\}$ shown in Figure 1. The symbol a_1 performs the circular shift in A_1 , and the identity in A_2 . Symmetrically, the symbol a_2 performs the identity in A_1 , and the circular shift in A_2 .

Figure 1: DFAs A_1 and A_2 with all valid states reachable in $\mathcal{D}(N)$.

Construct the NFA N recognizing the language $L(A_1)L(A_2)$ from the DFAs A_1 and A_2 by adding the transitions (f_1, a_2, s_2) and $(f_1 - 1, a_1, s_2)$, by making the state f_1 non-final and state s_2 non-initial. The NFA N is shown in Figure 2.

Figure 2: The NFA N recognizing the language $L(A_1)L(A_2)$.

Let us show that each valid state (j, S) is reachable in the subset automaton $\mathcal{D}(N)$. The proof is by induction on $|S|$. The basis, with $|S| = 0$, holds true since each state (j, \emptyset) with $j \leq n_1 - 1$ is reached from the initial state (s_1, \emptyset) by a_1^{j-1} . Let $|S| \geq 1$. There are three cases to consider.

Case 1: $j = f_1$. Then $s_2 \in S$ since (f_1, S) is valid. Since a_1 performs the circular shift in A_1 , and the identity in A_2 , we have $(n_1 - 1, S \setminus \{s_2\}) \xrightarrow{a_1} (f_1, \{s_2\} \cup (S \setminus \{s_2\})) = (f_1, S)$, where the leftmost state is reachable by induction.

Case 2: $j = s_1$. Let $m = \min S$. Then $s_2 \in a_2^{m-1}(S)$, and $|a_2^{m-1}(S)| = |S|$ since a_2 performs a permutation on the state set of A_2 . Since a_1 performs the identity on the state set of A_2 , we have

$$(f_1, a_2^{m-1}(S)) \xrightarrow{a_1} (s_1, a_2^{m-1}(S)) \xrightarrow{a_2^{m-1}} (s_1, S),$$

where the leftmost state is reachable as shown in Case 1.

Case 3: $2 \leq j \leq n_1 - 1$. Then we have $(s_1, S) \xrightarrow{a_1^{j-1}} (j, S)$, where the left state is considered in Case 2.

Thus, the two simple symbols a_1 and a_2 guarantee the reachability of all valid states in the subset automaton $\mathcal{D}(N)$. However, since both these symbols perform permutations on the state set Q_2 of A_2 , we have $Q_2 \cdot a_1 = Q_2 \cdot a_2 = Q_2$. It follows that in $\mathcal{D}(N)$, all states (i, Q_2) are equivalent to the final sink state.

To guarantee distinguishability, we add one more input symbol b which performs the contractions $s_1 \rightarrow 2$ and $s_2 \rightarrow 2$, and denote the resulting automata A'_1 and A'_2 , respectively. The NFA N' recognizing $L(A'_1)L(A'_2)$ is shown in Figure 3.

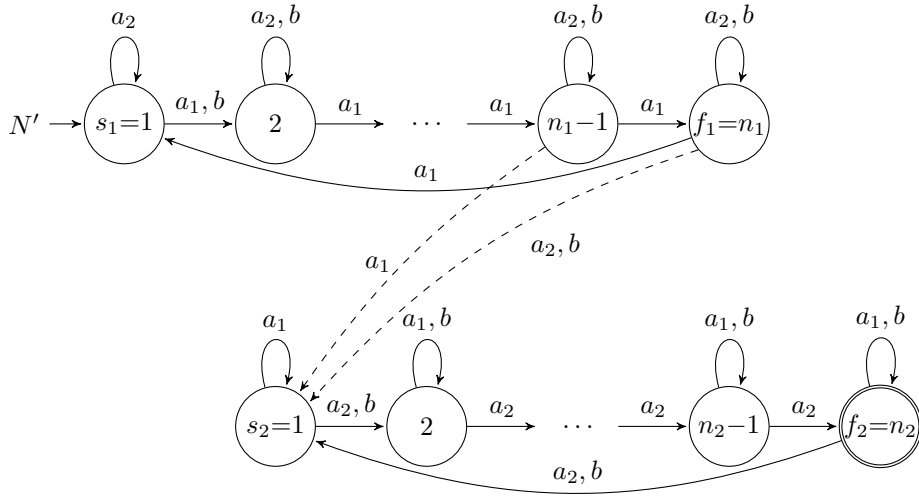


Figure 3: The NFA N recognizing the language $L(A'_1)L(A'_2)$.

As shown above, all valid states (j, S) are reachable in the corresponding subset automaton $\mathcal{D}(N')$. To get distinguishability, let us show that each singleton set is co-reachable in N' . In the reversed automaton $(N')^R$, the initial set is $\{f_2\}$, and

$$\{f_2\} \xrightarrow{a_2} \{n_2 - 1\} \xrightarrow{a_2} \{n_2 - 2\} \xrightarrow{a_2} \dots \xrightarrow{a_2} \{2\} \xrightarrow{a_2} \{s_2\}.$$

Next, since $n_1 \geq 3$, we have

$$\{s_2\} \xrightarrow{b} \{f_1\} \xrightarrow{a_1} \{n_1 - 1\} \xrightarrow{a_1} \dots \xrightarrow{a_1} \{2\} \xrightarrow{a_1} \{s_1\};$$

notice that we need $n_1 \geq 3$ to get $\{s_2\} \xrightarrow{b} \{f_1\}$, in the case of $n_1 = 2$ we would have $\{s_2\} \xrightarrow{b} \{f_1, s_1\}$. Hence each singleton set is co-reachable in N' . By Corollary 2, all states of the subset automaton $\mathcal{D}(N')$ are pairwise distinguishable. ■

We use the ideas from the above example to describe witnesses for multiple concatenation over a $(k + 1)$ -letter alphabet. To this aim, let $k \geq 2$ and $n_i \geq 3$ for $i = 1, 2, \dots, k$. Let $\Sigma = \{b, a_1, a_2, \dots, a_k\}$ be an alphabet consisting of $k + 1$ symbols. Define an n_i -state DFA $A_i = (Q_i, \Sigma, \cdot, s_i, \{f_i\})$, where

- $Q_i = \{1, 2, \dots, n_i\}$,
- $s_i = 1$,
- $f_i = n_i$,
- $a_i: (1, 2, \dots, n_i)$, $a_j: (1)$ if $j \neq i$, $b: (1 \rightarrow 2)$,

that is, the symbol a_i performs the circular shift on Q_i , each symbol a_j with $j \neq i$ performs the identity, and the symbol b performs a contraction. The DFA A_i is shown in Figure 4; here $\Sigma \setminus \{a_i\}$ on a loop means that there is a loop in the corresponding state on each symbol in $\Sigma \setminus \{a_i\}$, and the same for $\Sigma \setminus \{a_i, b\}$.

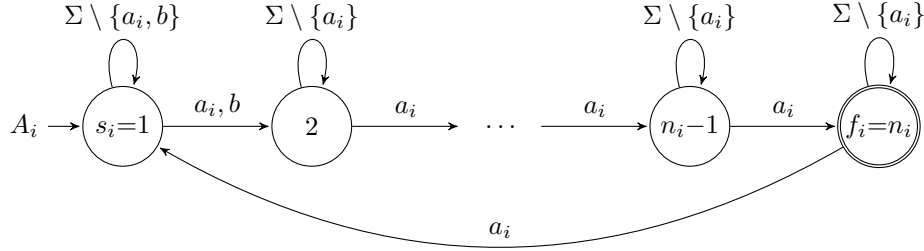


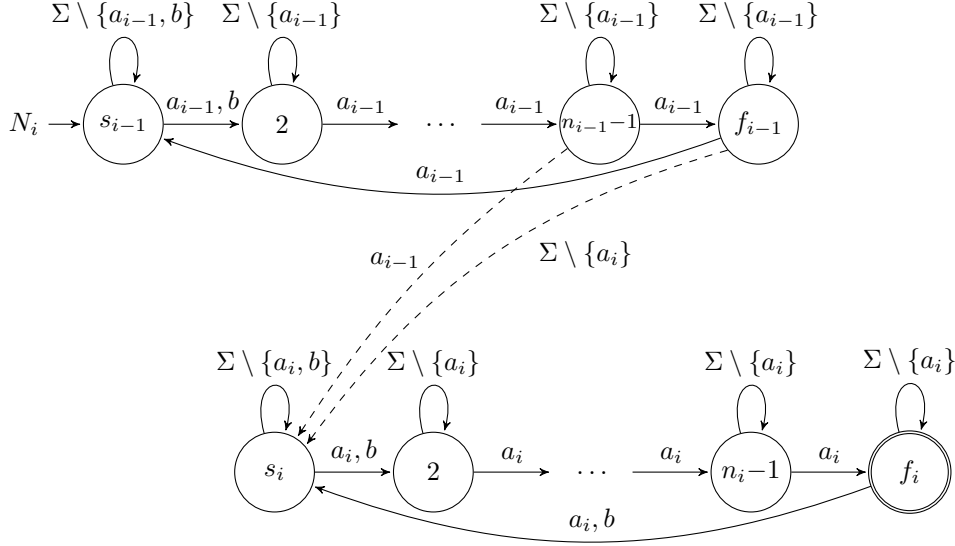
Figure 4: The witness DFA A_i over the $(k + 1)$ -letter alphabet $\{b, a_1, a_2, \dots, a_k\}$.

First, let us consider the concatenation $L(A_{i-1})L(A_i)$ where $2 \leq i \leq k$. Construct an NFA N_i for this concatenation from DFAs A_{i-1} and A_i as shown in Figure 5, that is, by adding the transitions $(f_{i-1}-1, a_{i-1}, s_i)$ and (f_{i-1}, σ, s_i) with $\sigma \in \Sigma \setminus \{a_{i-1}\}$, by making the state f_{i-1} non-final, and the state s_i non-initial.

The next observation is crucial in what follows. It shows that in the subset automaton $\mathcal{D}(N_i)$, each state (s_{i-1}, S) with $S \subseteq Q_i$ and $S \neq \emptyset$ is reachable from $(s_{i-1}, \{s_i\})$. Moreover, while reaching (s_{i-1}, S) with $f_i \notin S$, the state f_i is never visited. This is a very important property since, later, we do not wish to influence the $(i + 1)$ st component of a valid state while setting its i th component.

Lemma 11. *Let $2 \leq i \leq k$ and N_i be the NFA for the language $L(A_{i-1})L(A_i)$ described above. For every non-empty subset $S \subseteq Q_i$, there exists a string w_S over the alphabet $\{a_{i-1}, a_i\}$ such that in the subset automaton $\mathcal{D}(N_i)$, we have*

- (i) $(s_{i-1}, \{s_i\}) \xrightarrow{w_S} (s_{i-1}, S)$;
- (ii) if $f_i \notin S$, u is a prefix of w_S , and $(s_{i-1}, \{s_i\}) \xrightarrow{u} (q, T)$, then $f_i \notin T$.

Figure 5: The NFA N_i recognizing the language $L(A_{i-1})L(A - i)$.

Proof. The proof of both (i) and (ii) is by induction on $|S|$. The basis, with $|S| = 1$, holds true since we have

$$(s_{i-1}, \{s_i\}) \xrightarrow{a_i} (s_{i-1}, \{2\}) \xrightarrow{a_i} \dots \xrightarrow{a_i} (s_{i-1}, \{n_i - 1\}) \xrightarrow{a_i} (s_{i-1}, \{f_i\}),$$

so, for each $j \in Q_i$, the state $(s_{i-1}, \{j\})$ is reached from $(s_{i-1}, \{s_i\})$ by a_i^{j-1} . Moreover, if $j \neq f_i$, then f_i is not visited while reading a_i^{j-1} .

Let $|S| \geq 2$. Let $m = \min S$ and $S' = a_i^{m-1}(S \setminus \{m\})$. Then $|S'| = |S| - 1$. By reading n_{i-1} times the symbol a_{i-1} and then the string a_i^{m-1} we get

$$(s_{i-1}, S') \xrightarrow{a_{i-1}^{n_{i-1}}} (s_{i-1}, \{s_i\} \cup S') \xrightarrow{a_i^{m-1}} (s_{i-1}, \{m\} \cup (S \setminus \{m\})) = (s_{i-1}, S),$$

where the leftmost state is reached from $(s_{i-1}, \{s_i\})$ by the string $w_{S'}$ by induction, so $w_S = w_{S'} a_{i-1}^{n_{i-1}} a_i^{m-1}$. Moreover, if $f_i \notin S$, then $S' \subseteq [2, f_i - m]$, so $f_i \notin S'$. By induction, the state f_i has not been visited while reading $w_{S'}$ to reach (s_{i-1}, S') from $(s_{i-1}, \{s_i\})$. Since in A_i , the symbols a_{i-1} and a_i perform the identity and circular shift, respectively, the state f_i is not visited either while reading the string $a_{i-1}^{n_{i-1}} a_i^{m-1}$ to reach (s_{i-1}, S) from (s_{i-1}, S') . \square

Now, construct the NFA N recognizing the concatenation $L(A_1)L(A_2)\dots L(A_k)$ from DFAs A_1, A_2, \dots, A_k as follows: First, for each $i = 1, 2, \dots, k-1$, add the transitions (f_{i-1}, a_i, s_{i+1}) and (f_i, σ, s_{i+1}) with $\sigma \in \Sigma \setminus \{a_i\}$. Then, make states f_1, f_2, \dots, f_{k-1} non-final, and states s_2, s_3, \dots, s_k non-initial; see Figure 6 for an illustration.

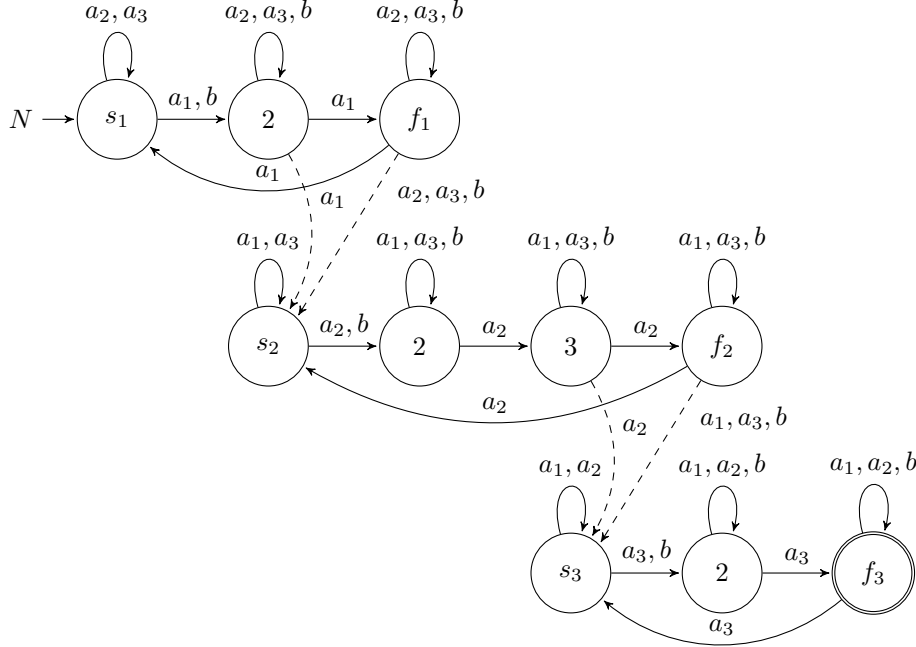


Figure 6: The NFA N for $L(A_1)L(A_2)L(A_3)$ with $n_1 = 3$, $n_2 = 4$, and $n_3 = 3$.

Theorem 12. *Let $k \geq 2$ and $n_i \geq 3$ for $i = 1, 2, \dots, k$. Let A_i be the n_i -state DFA from Figure 4. Let N be the NFA for $L(A_1)L(A_2) \cdots L(A_k)$ described above. Then all valid states are reachable and pairwise distinguishable in the subset automaton $\mathcal{D}(N)$.*

Proof. We first prove reachability. Let $q = (j, S_2, S_3, \dots, S_k)$ be a valid state. If $S_2 = \emptyset$, then the state $q = (j, \emptyset, \emptyset, \dots, \emptyset)$ is reached from the initial state $(s_1, \emptyset, \emptyset, \dots, \emptyset)$ by the string a_1^{j-1} . Next, let $\ell = \max\{i \geq 2 \mid S_i \neq \emptyset\}$. Then $q = (j, S_2, S_3, \dots, S_\ell, \emptyset, \emptyset, \dots, \emptyset)$ where $2 \leq \ell \leq k$, $S_i \subseteq Q_i$ and $S_i \neq \emptyset$ for $i = 2, 3, \dots, \ell$. Since each a_i performs the circular shift in A_i and the identity in A_j with $j \neq i$, the string $a_1^{n_1} a_2^{n_2} \cdots a_{\ell-1}^{n_{\ell-1}}$ sends the initial state $(s_1, \emptyset, \emptyset, \dots, \emptyset)$ to

$$(s_1, \{s_2\}, \{s_3\}, \dots, \{s_{\ell-1}\}, \{s_\ell\}, \emptyset, \emptyset, \dots, \emptyset).$$

Now, we are going to set the corresponding components to sets S_i , starting with S_ℓ , continuing with $S_{\ell-1}, S_{\ell-2}, \dots$, and ending with S_3 and S_2 . By Lemma 11 applied to the NFA N_ℓ recognizing the language $L(A_{\ell-1})L(A_\ell)$, there is a string w_{S_ℓ} over $\{a_{\ell-1}, a_\ell\}$ which sends $(s_{\ell-1}, \{s_\ell\})$ to $(s_{\ell-1}, S_\ell)$ in the subset automaton $\mathcal{D}(N_\ell)$. Moreover, since q is valid, we have $f_\ell \notin S_\ell$, which means that the state f_ℓ is not visited while reading w_{S_ℓ} . Since both $a_{\ell-1}$ and a_ℓ perform identities on $Q_1, Q_2, \dots, Q_{\ell-2}$, in the subset automaton $\mathcal{D}(N)$ we have

$$(s_1, \{s_2\}, \dots, \{s_{\ell-1}\}, \{s_\ell\}, \emptyset, \dots, \emptyset) \xrightarrow{w_{S_\ell}} (s_1, \{s_2\}, \dots, \{s_{\ell-1}\}, S_\ell, \emptyset, \dots, \emptyset).$$

Next, Lemma 11 applied to $N_{\ell-1}$ gives a string $w_{S_{\ell-1}}$ over $\{a_{\ell-2}, a_{\ell-1}\}$ which sends $(s_{\ell-2}, \{s_{\ell-1}\})$ to $(s_{\ell-2}, S_{\ell-1})$ in $\mathcal{D}(N_{\ell-1})$, and moreover if $f_{\ell-1} \notin S_{\ell-1}$, then $f_{\ell-1}$ is not visited while reading this string. Since both symbols $a_{\ell-2}$ and $a_{\ell-1}$ perform identities on $Q_1, Q_2, \dots, Q_{\ell-3}$, as well as on Q_ℓ , in $\mathcal{D}(N)$ we have

$$(s_1, \{s_2\}, \dots, \{s_{\ell-2}\}, \{s_{\ell-1}\}, S_\ell, \emptyset, \dots, \emptyset) \xrightarrow{w_{S_{\ell-1}}} (s_1, \{s_2\}, \dots, \{s_{\ell-2}\}, S_{\ell-1}, S_\ell, \emptyset, \dots, \emptyset).$$

Now, for $i = 2, 3, \dots, \ell-2$, let w_{S_i} be the string over $\{a_{i-1}, a_i\}$ given by Lemma 11 that sends $(s_{i-1}, \{s_i\})$ to (s_{i-1}, S_i) in the NFA N_i for $L(A_{i-1})L(A_i)$. Moreover, $f_i \notin S_i$ implies that the state f_i is never visited while reading w_{S_i} , which in turn implies that s_{i+1} is never added to the $(i+1)$ th component in such a case. If $f_i \in S_i$ and $i \leq k-1$, then the state s_{i+1} is included in S_{i+1} since the state q is valid, and s_{i+1} is sent to itself by both a_{i-1} and a_i . Next, there is a loop on both symbols a_{i-1} and a_i in the states s_1, s_2, \dots, s_{i-2} , as well as in all states of automata $A_{i+1}, A_{i+2}, \dots, A_\ell$. Set $W = w_{S_{\ell-2}} w_{S_{\ell-3}} \dots w_{S_3} w_{S_2}$. Then in $\mathcal{D}(N)$ we have

$$(s_1, \{s_2\}, \dots, \{s_{\ell-2}\}, S_{\ell-1}, S_\ell, \emptyset, \dots, \emptyset) \xrightarrow{W} (s_1, S_2, \dots, S_{\ell-2}, S_{\ell-1}, S_\ell, \emptyset, \dots, \emptyset),$$

and the resulting state is sent to the state q by the string a_1^{j-1} . Hence the valid state $q = (j, S_2, S_3, \dots, S_\ell, \emptyset, \emptyset, \dots, \emptyset)$ is reached from the initial state $(s_1, \emptyset, \emptyset, \dots, \emptyset)$ by the string $a_1^{n_1} a_2^{n_2} \dots a_{\ell-1}^{n_{\ell-1}} w_{S_\ell} w_{S_{\ell-1}} \dots w_{S_3} w_{S_2} a_1^{j-1}$.

To get distinguishability, let us show that each singleton set is co-reachable in N . First, for an example, consider the NFA from Figure 6. In its reversed automaton, the initial set is $\{f_3\}$, and we have

$$\{f_3\} \xrightarrow{a_3} \{2\} \xrightarrow{a_3} \{s_3\} \xrightarrow{b} \{f_2\} \xrightarrow{a_2} \{3\} \xrightarrow{a_2} \{2\} \xrightarrow{a_2} \{s_2\} \xrightarrow{b} \{f_1\} \xrightarrow{a_1} \{2\} \xrightarrow{a_1} \{s_1\}.$$

In the general case, the initial set of N^R is $\{f_k\}$. Next, for each $i = 1, 2, \dots, k$, each singleton set $\{j\}$ such that $j \in Q_i$ is reached from $\{f_i\}$ via a string in a_i^* . Finally, for each $i = 2, 3, \dots, k$, the singleton set $\{f_{i-1}\}$ is reached from $\{s_i\}$ by b since $n_{i-1} \geq 3$. Thus, for every state q of N , the singleton set $\{q\}$ is co-reachable in the NFA N . By Corollary 2, all states of the subset automaton $\mathcal{D}(N)$ are pairwise distinguishable. \square

Notice that all automata in the previous theorem, as well as witness automata from [1], are required to have at least three states. We conclude this section by describing the witnesses for multiple concatenation also in the case where some of given automata have two states. The idea is to use symbols a_k and b to guarantee co-reachability of singleton sets in such a way that they perform either the identity or $(1 \rightarrow 2 \rightarrow \dots \rightarrow n_i)$ in every second automaton. However, then we should be careful with reachability. To this aim, let $k \geq 2$, $n_i \geq 2$ for $i = 1, 2, \dots, k$, and $\Sigma = \{b, a_1, a_2, \dots, a_k\}$. Let

$$I = \{i \mid 1 \leq i \leq k-1 \text{ and } i \bmod 2 = k \bmod 2\}$$

$$J = \{i \mid 1 \leq i \leq k-1 \text{ and } i \bmod 2 \neq k \bmod 2\},$$

that is, the set I contains the indexes that have the same parity as k , and the set J the others.

Consider the n_i -state DFAs $A_i = (Q_i, \Sigma, \cdot, s_i, \{f_i\})$, see Figure 7, where we have $Q_i = \{1, 2, \dots, n_i\}$, $s_i = 1$, $f_i = n_i$, and the transitions are as follows:
 if $i \in I$, then $a_i: (1, 2, \dots, n_i)$, $a_k: (1 \rightarrow 2 \rightarrow \dots \rightarrow n_i)$, and $\sigma: (1)$ if $\sigma \in \Sigma \setminus \{a_i, a_k\}$,
 if $i \in J$, then $a_i: (1, 2, \dots, n_i)$, $b: (1 \rightarrow 2 \rightarrow \dots \rightarrow n_i)$, and $\sigma: (1)$ if $\sigma \in \Sigma \setminus \{a_i, b\}$,
 if $i = k$, then $b: (1, 2, \dots, n_k)$, $a_k: (1 \rightarrow 2 \rightarrow \dots \rightarrow n_k)$, and $\sigma: (1)$ if $\sigma \in \Sigma \setminus \{a_k, b\}$,
 that is,

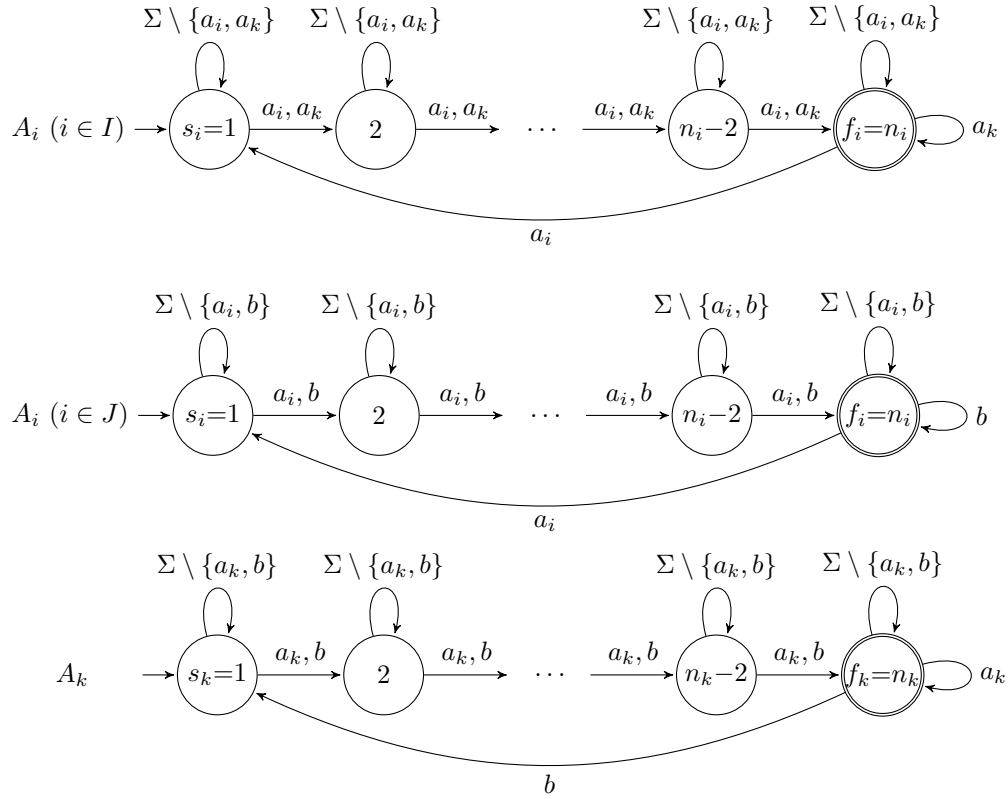


Figure 7: The DFAs A_i with $i \in I$ (top), A_i with $i \in J$ (middle), and A_k (bottom).

- each a_i with $1 \leq i \leq k-1$ performs the circular shift on Q_i , and the identity on Q_j with $j \neq i$;
- a_k performs the transformation $(1 \rightarrow 2 \rightarrow 3 \rightarrow \dots \rightarrow n_i)$ on Q_i with $i \in I$ or $i = k$, and the identity on Q_i with $i \in J$,
- b performs the transformation $(1 \rightarrow 2 \rightarrow 3 \rightarrow \dots \rightarrow n_i)$ on Q_i with $i \in J$, the circular shift on Q_k , and the identity on Q_i with $i \in I$.

Construct an NFA N for the language $L(A_1)L(A_2)\cdots L(A_k)$ from the DFAs A_1, A_2, \dots, A_k as follows (see Figure 8 for an illustration):

For each $i = 1, 2, \dots, k-1$, add the transitions (f_{i-1}, a_i, s_{i+1}) and (f_i, σ, s_{i+1}) for each $\sigma \in \Sigma \setminus \{a_i\}$, and moreover, if $i \in I$, then add the transition (f_{i-1}, a_k, s_{i+1}) , and if $i \in J$, then add the transition (f_{i-1}, b, s_{i+1}) . The initial state of N is s_1 , and its unique final state is f_k .

Theorem 13. *Let $k \geq 2$ and $n_i \geq 2$ for $i = 1, 2, \dots, k$. Let A_1, A_2, \dots, A_k be the DFAs shown in Figure 7, and N be the NFA for $L(A_1)L(A_2)\cdots L(A_k)$ described above. Then all valid states are reachable and pairwise distinguishable in $\mathcal{D}(N)$.*

Proof. First, notice that Lemma 11 still holds for automata A_1, A_2, \dots, A_{k-1} since the transitions on a_1, a_2, \dots, a_{k-1} are the same. Thus, for each non-empty subset S of Q_i with $i \leq k-1$, let w_S be the string over $\{a_{i-1}, a_i\}$ given By Lemma 11.

Let $(\{j\}, S_2, S_3, \dots, S_k)$ be a valid state. If $S_k = \emptyset$, then $(j, S_2, S_3, \dots, S_{k-1}, \emptyset)$ is reachable as shown in the proof of Theorem 12.

Now, let $S_k \neq \emptyset$. Then the state $(s_1, \{s_2\}, \{s_3\}, \dots, \{s_k\})$ is reached from the initial state by $a_1^{n_1} a_2^{n_2} \cdots a_{k-1}^{n_{k-1}}$. Next, notice that Lemma 11 still holds for N_k even if a_k fixes f_k instead of sending it to s_k since the out-transition in f_k on a_k is not used in the proof of the lemma. Hence, there is a string $w(S_k)$ over $\{a_{k-1}, a_k\}$ which sends the state $(s_{k-1}, \{s_k\})$ to (s_{k-1}, S_k) in the subset automaton $\mathcal{D}(N_k)$. However, each a_k sends each state s_i with $i \in I$ to $s_i + 1$, and we must then read the string $u_i = (a_i)^{n_i-1}$ to send $s_i + 1$ back to s_i while fixing the states in all the remaining components. Let $u = \prod_{i \in I} u_i$. Now, let $w'(S_k)$ be the string obtained from $w(S_k)$ by inserting u after each a_k . Since before reading each a_k in w_{S_k} we have s_{k-1} in the $(k-1)$ st component, the state $(s_1, \{s_2\}, \dots, \{s_{k-1}\}, \{s_k\})$ is sent to $(s_1, \{s_2\}, \dots, \{s_{k-1}\}, S_k)$ by w'_{S_k} , and then to $(j, S_2, S_3, \dots, S_{k-1}, S_k)$ by $w_{S_{k-1}} w_{S_{k-2}} \cdots w_{S_3} w_{S_2} a_1^{j-1}$.

To prove distinguishability, let us show that all singleton sets are co-reachable in the NFA N . First, as an example, consider the NFA N from Figure 8, and notice that in the reversed automaton N^R , we have

$$\begin{aligned} \{f_5\} &\xrightarrow{b} \{s_5\} \xrightarrow{a_5} \{f_4\} \xrightarrow{a_4} \{s_4\} \xrightarrow{b} \{f_3\} \xrightarrow{a_3} \{2\} \xrightarrow{a_3} \{s_3\} \\ &\xrightarrow{a_5} \{f_2\} \xrightarrow{a_2} \{s_2\} \xrightarrow{b} \{f_1\} \xrightarrow{a_1} \{2\} \xrightarrow{a_1} \{s_1\}. \end{aligned}$$

In the general case, the initial set of the reversed automaton N^R is $\{f_k\}$, and each set $\{q\}$ with $q \in Q_k$ is reached from $\{f_k\}$ by a string in b^* . Next each $\{f_i\}$ with $i \in J$ is reached from $\{s_{i+1}\}$ by a_k , while each $\{f_i\}$ with $i \in I$ is reached from $\{s_{i+1}\}$ by b . Finally, each $\{q\}$ with $q \in Q_i$, where $1 \leq i \leq k-1$, is reached from $\{f_i\}$ by a string in a_i^* . It follows that all singleton sets are co-reachable in N . By Corollary 2, all states of $\mathcal{D}(N)$ are pairwise distinguishable. \square

5. Matching Lower Bound: k -letter Alphabet

The aim of this section is to describe witnesses for multiple concatenation over a k -letter alphabet. Let us start with the following example.

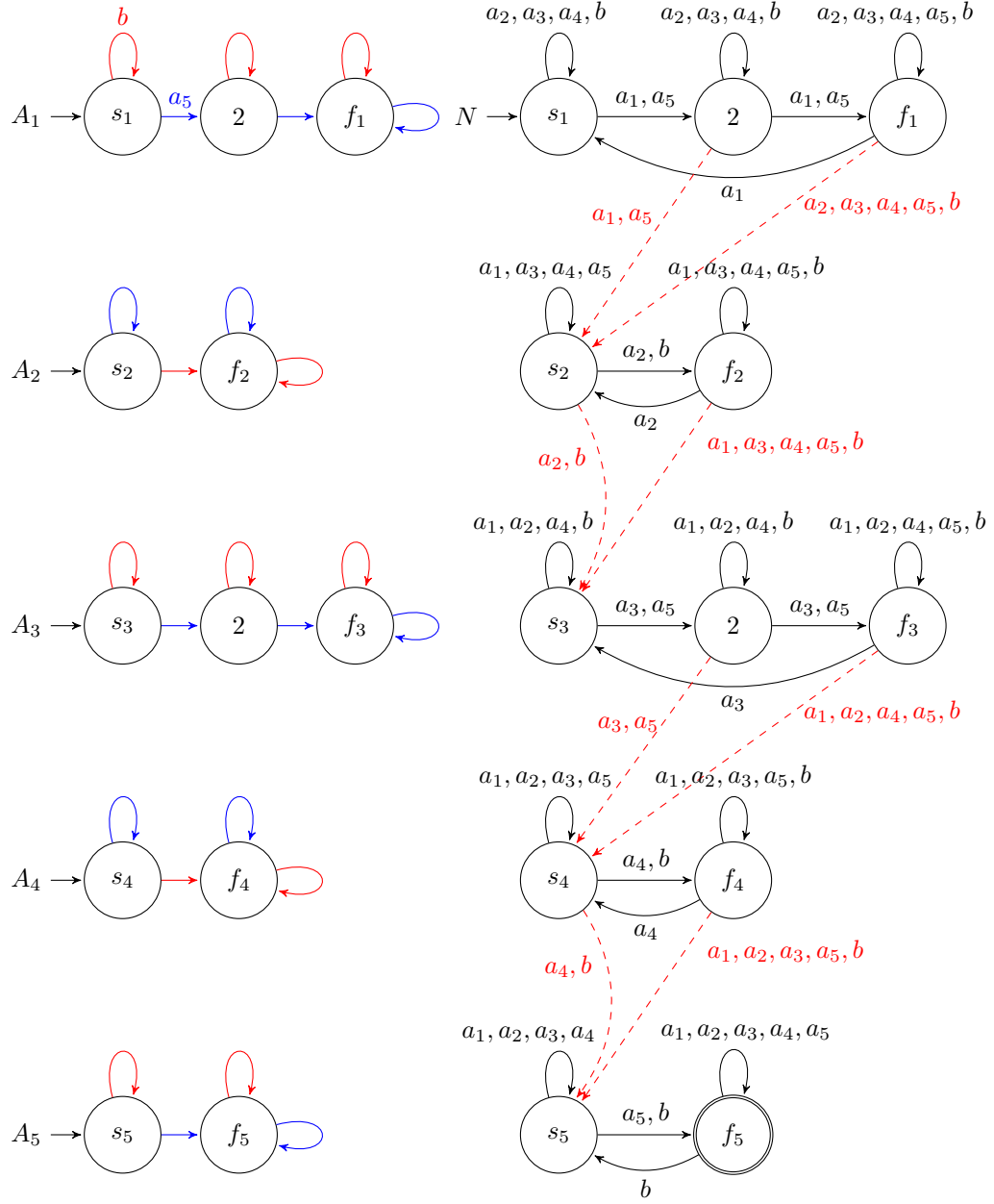


Figure 8: The DFAs A_1, A_2, A_3, A_4, A_5 : transitions on a_5 and b (left) and the NFA N for $L(A_1)L(A_2)L(A_3)L(A_4)L(A_5)$ (right) with $n_1 = n_3 = 3$ and $n_2 = n_4 = n_5 = 2$.

Example 14. Let $n_1, n_2 \geq 1$ and A and B be the binary DFAs shown in Figure 9. Let us show that the languages $L(A)$ and $L(B)$ are witnesses for concatenation of two regular languages.

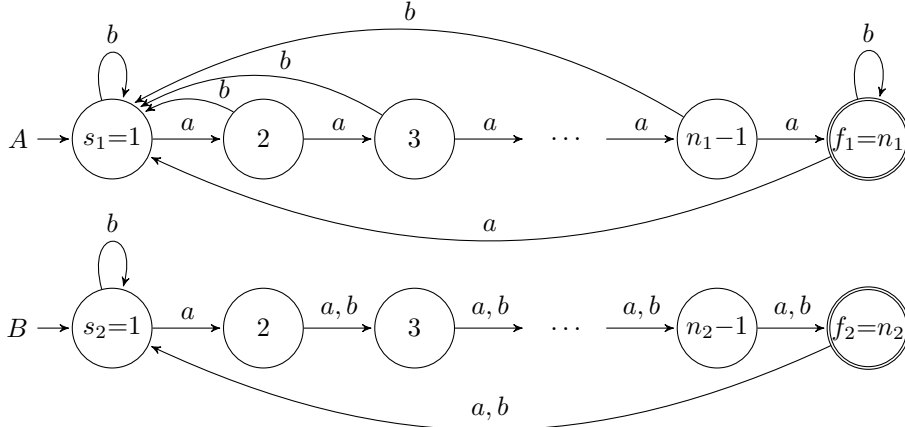


Figure 9: The binary witnesses for concatenation; $n_1, n_2 \geq 1$.

First, let $n_2 = 1$. Then $L(B) = \{a, b\}^*$ and the concatenation $L(A)\{a, b\}^*$ is recognized by the minimal n_1 -state DFA obtained from A by replacing the transition (f_1, a, s_1) with the transition (f_1, a, f_1) . An upper bound is n_1 by Proposition 8.

Now, let $n_1 = 1$ and $n_2 \geq 2$. Then $s_1 = f_1$. Construct an NFA N for $L(A)L(B)$ from the DFAs A and B by adding the transitions (f_1, a, s_2) and (f_1, b, s_2) , and by making the state s_1 non-final. Let us show that all valid states (f_1, S) are reachable in $\mathcal{D}(N)$. Since (f_1, S) is valid, we have $s_2 \in S$. The proof is by induction on $|S|$. The basis, $|S| = 1$, that is, $S = \{s_2\}$, holds true since $(f_1, \{s_2\})$ is the initial state. Let $|S| \geq 2$ and $s_2 \in S$. Let $m = \min(S \setminus \{s_2\})$ and $S' = S \setminus \{s_2, m\}$. Then $ab^{m-2}(S') \subseteq [2, n_2 - m + 1]$ and

$$(f_1, \{s_2\} \cup ab^{m-2}(S')) \xrightarrow{a} (f_1, \{s_2, 2\} \cup b^{m-2}(S')) \xrightarrow{b^{m-2}} (f_1, \{s_2, m\} \cup S') = (f_1, S),$$

where the leftmost valid state is reachable by induction. This proves the reachability of 2^{n_2-1} valid states. All these states are pairwise distinguishable by Lemma 1 since all singletons $\{q\}$, where q is a state of B , are co-reachable in N . By Proposition 8, an upper bound is $V_2 = 2^{n_2-1}$.

Finally, let $n_1, n_2 \geq 2$. Construct an NFA N for $L(A)L(B)$ from the DFAs A and B by adding the transitions (f_1-1, a, s_2) and (f_1-1, b, s_2) , by making the state f_1 non-final and the state s_2 non-initial. Let us show that in the subset automaton $\mathcal{D}(N)$, each valid state (j, S) is reachable. The proof is by induction on $|S|$. The basis, with $|S| = 0$, holds true since each valid state (j, \emptyset) is reached from the initial state is (s_1, \emptyset) by a^{j-1} . Let $|S| \geq 1$. There are three cases to consider.

Case 1: $j = f_1$. Then $s_2 \in S$ since (f_1, S) is valid. We have

$$(f_1 - 1, a(S \setminus \{s_2\})) \xrightarrow{a} (f_1, \{s_2\}) \cup (S \setminus \{s_2\}) = (f_1, S)$$

where the leftmost valid state is reachable by induction.

Case 2: $j = s_1$.

Case 2.a: $2 \in S$. Then $s_2 \in a(S)$ and (s_1, S) is reached from $(f_1, a(S))$ by a , where the latter valid state is considered in Case 1.

Case 2.b: $2 \notin S$ and $S = \{s_2\}$. Then we have $(f_1, \{s_2\}) \xrightarrow{a} (s_1, \{2\}) \xrightarrow{b^{n_2}} (s_1, \{s_2\})$, where the leftmost state is considered in Case 1.

Case 2.c: $2 \notin S$ and $S \neq \{s_2\}$. Let $m = \min(S \setminus \{s_2\})$ and $S' = S \cap \{s_2\}$. Then $2 \in b^{m-2}(S \setminus \{s_2\})$ and (s_1, S) is reached from $(s_1, S' \cup b^{m-2}(S \setminus \{s_2\}))$ by b^{m-2} where the latter state is considered in Case 2.a.

Case 3: $2 \leq j \leq n_1 - 1$. Then (j, S) is reached from $(s_1, a^{j-1}(S))$ by a^{j-1} , and the latter set is considered in Case 2.

This proves the reachability of $(n_1 - 1)2^{n_2} + 2^{n_2-1}$ states. To get distinguishability, let (i, S) and (j, T) be two distinct valid states. There are two cases to consider.

Case 1: $S \neq T$. The two states are distinguishable by Lemma 1 since all singletons $\{q\}$, where q is a state of B , are co-reachable in N .

Case 2: $S = T$ and $i < j$. First, let $S = \emptyset$. Since $n_1 \geq 2$, the string a^{n_1-j} sends the two states to states that differ in s_2 . The resulting states are distinguishable as shown in Case 1. Now, let $S \neq \emptyset$. Then the two states are sent to $(s_1, \{s_2\})$ and $(f_1, \{s_2\})$ by $a^{n_1-j}b^{n_2}$. Let us show that the resulting states are sent to states that differ in s_2 by a^{n_1} if $s_2a^{n_1} \neq s_2$, and by $a^{n_1-1}ba^{n_1-1}$ otherwise.

First, notice that both strings a^{n_1} and $a^{n_1-1}ba^{n_1-1}$ send the state f_1 to itself in A . It follows that $(f_1, \{s_2\})$ is sent to a state containing s_2 in its second component by both these strings.

Now, let $s_2a^{n_1} \neq s_2$. Then we have

$$(s_1, \{s_2\}) \xrightarrow{a^{n_1-1}} (f_1, \{s_2, s_2a^{n_1-1}\}) \xrightarrow{a} (s_1, \{s_2a, s_2a^{n_1}\}),$$

where $s_2a \neq s_2$ since $n_2 \geq 2$. Thus, in this case, the string a^{n_1} sends the state $(s_1, \{s_2\})$ to a state which does not have s_2 in its second component.

Finally, let $s_2a^{n_1} = s_2$. Then $s_2a^{n_1-1} = f_2$ and since $s_2b = f_2b = s_2$, we have

$$(s_1, \{s_2\}) \xrightarrow{a^{n_1-1}} (f_1, \{s_2, f_2\}) \xrightarrow{b} (f_1, \{s_2\}) \xrightarrow{a^{n_1-1}} (f_1 - 1, \{f_2\}),$$

where $f_2 \neq s_2$ since $n_2 \geq 2$. Hence, this time the string $a^{n_1-1}ba^{n_1-1}$ sends $(s_1, \{s_2\})$ to a state which does not contain s_2 in its second component.

This proves distinguishability, and concludes our proof since by Theorem 6, a (known) upper bound is $(n_1 - 1)U_2 + V_2 = (n_1 - 1)2^{n_2} + 2^{n_2-1}$ in this case. \blacksquare

Hence the above example provides a two-letter witnesses for the concatenation of two regular languages (even in the case then automata may have one or two states). Therefore, in what follows we assume that $k \geq 3$.

We use our previous results to describe witnesses for the concatenation of k languages over the k -letter alphabet $\{b, a_1, a_2, \dots, a_{k-1}\}$. The idea is as follows. The transitions on input symbols a_1, a_2, \dots, a_{k-1} in automata A_1, A_2, \dots, A_{k-1} are the same as in our $(k+1)$ -letter witnesses from Theorem 12, while A_{k-1} and A_k over $\{a_{k-1}, b\}$ are the same as automata A and B in Example 14. The input symbol b performs the transformation $(\{2, 3, \dots, n_i - 1\} \rightarrow s_i)$ in each A_i except for A_k , and it is used to get reachability as well as distinguishability.

To this aim, let $k \geq 3$ and $\Sigma = \{b, a_1, a_2, \dots, a_{k-1}\}$ be a k -letter alphabet. Let $n_1, n_k \geq 2$ and $n_i \geq 3$ for $i = 2, 3, \dots, k-1$. For $i = 1, 2, \dots, k$, define an n_i -state DFA $A_i = (Q_i, \Sigma, \cdot, s_i, \{f_i\})$, see Figure 10, where $Q_i = \{1, 2, \dots, n_i\}$, $s_i = 1$, $f_i = n_i$, and the transitions are as follows:

- if $i \leq k-1$, then
 $a_i: (1, 2, \dots, n_i)$, $b: (\{2, 3, \dots, n_i - 1\} \rightarrow s_i)$, and $\sigma: (1)$ if $\sigma \in \Sigma \setminus \{a_i, b\}$,
- if $i = k$, then
 $a_{k-1}: (1, 2, \dots, n_k)$, $b: (2 \rightarrow 3 \rightarrow \dots \rightarrow n_k \rightarrow 1)$, and $\sigma: (1)$ if $\sigma \in \Sigma \setminus \{a_{k-1}, b\}$.

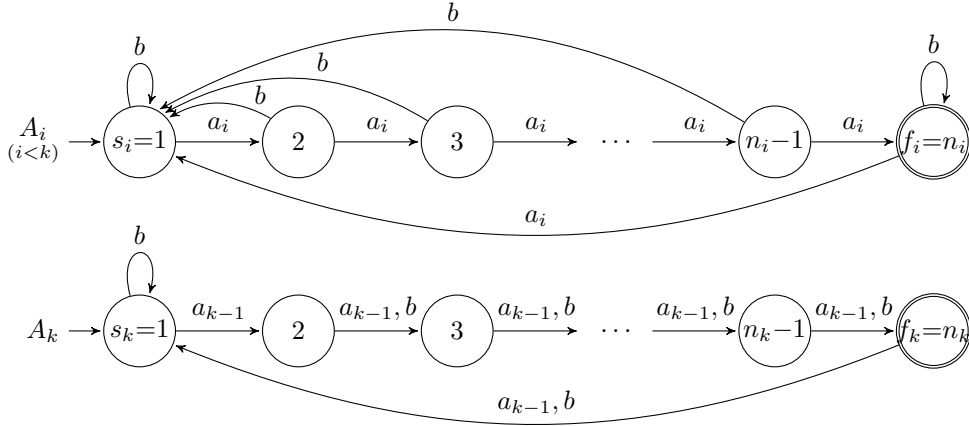


Figure 10: The DFA A_i with $i < k$ (top): transitions on a_i and b , and the DFA A_k (bottom): transitions on a_{k-1} and b ; all the remaining symbols in both automata perform identities; $n_1, n_k \geq 2$ and $n_i \geq 3$ for $i = 2, 3, \dots, k-1$.

Construct an NFA N for $L(A_1)L(A_2)\dots L(A_k)$ from DFAs A_1, A_2, \dots, A_k by adding the transitions (f_i-1, a_i, s_{i+1}) , (f_i, a_j, s_{i+1}) for $j \neq i$, and (f_i, b, s_{i+1}) for $i = 1, 2, \dots, k-1$; the initial state of N is s_1 , and the final state is f_k . The next theorem shows that all valid states are reachable and pairwise distinguishable in $\mathcal{D}(N)$. The proof of reachability is based on our results concerning $(k+1)$ -letter witnesses as well as our binary witnesses from Example 14. The proof of distinguishability is not for free this time.

Theorem 15. *Let $k \geq 3$, $n_1, n_k \geq 2$, and $n_i \geq 3$ for $i = 2, 3, \dots, k-1$. Let A_1, A_2, \dots, A_k be DFAs shown in Figure 10 over the k -letter alphabet $\{b, a_1, a_2, \dots, a_{k-1}\}$. Let N be the NFA for $L(A_1)L(A_2) \cdots L(A_k)$ described above. Then all valid states are reachable and pairwise distinguishable in $\mathcal{D}(N)$.*

Proof. Consider a valid state $q = (j, S_2, \dots, S_{k-1}, S_k)$. First, let $S_k = \emptyset$. Since the transitions on a_1, a_2, \dots, a_{k-1} in A_1, A_2, \dots, A_{k-1} are the same as in automata in Theorem 12, the valid state $(j, S_2, \dots, S_{k-1}, \emptyset)$ is reachable exactly the same way as in the proof of this theorem.

Now let $S_k \neq \emptyset$. Notice that the transitions on a_{k-1} and b in DFAs A_{k-1} and A_k are the same as those on a and b in DFAs A and B in Example 14. As shown in this example, for each $S \subseteq Q_k$, there is a string w_S over $\{a_{k-1}, b\}$ which sends (s_{k-1}, \emptyset) to (s_{k-1}, S) in the subset automaton for $L(A_{k-1})L(A_k)$. Since we have a loop on both a_{k-1} and b in all states s_1, s_2, \dots, s_{k-2} , we reach $(s_1, \{s_2\}, \{s_3\}, \dots, \{s_{k-2}\}, \{s_{k-1}\}, S)$ from the initial state by $a_1^{n_1} a_2^{n_2} \cdots a_{k-2}^{n_{k-2}} w_S$. Next, let w_{S_i} be the string over $\{a_{i-2}, a_{i-1}\}$ given by Lemma 11 which sends $(s_{i-1}, \{s_i\})$ to (s_{i-1}, S_i) . Recall that $f_i \notin S_i$ implies that the state f_i is not visited while reading w_{S_i} . Moreover, a closer look at the proof of the lemma shows that if $f_i \in S_i$ then f_i is visited for the first time immediately after reading the last a_i in w_{S_i} . Now, let m be the number of occurrences of the symbol a_{k-1} in the string $w_{S_{k-1}}$. Then the state $(s_1, \{s_2\}, \{s_3\}, \dots, \{s_{k-2}\}, \{s_{k-1}\}, a_{k-1}^m(S_k))$ is reachable as shown above, and it is sent to $(s_1, \{s_2\}, \{s_3\}, \dots, \{s_{k-2}\}, S_{k-1}, S_k)$ by $w_{S_{k-1}}$. The resulting state is sent to q by the string $w_{S_{k-2}} w_{S_{k-3}} \cdots w_{S_3} w_{S_2} a_1^{j-1}$.

To get distinguishability, let $p = (S_1, S_2, S_3, \dots, S_k)$ and $q = (T_1, T_2, T_3, \dots, T_k)$ be two distinct valid states. If $S_k \neq T_k$, then p and q are distinguishable by Lemma 1 since each singleton subset of Q_k is co-reachable in N via a string in a_{k-1}^* .

Let $S_i \neq T_i$ for some i with $1 \leq i \leq k-1$, and $S_j = T_j$ for $j = i+1, i+2, \dots, k$. Let us show that there is a string that sends p and q to two states which differ in s_{i+1} .

Without loss of generality, we have $s \in S_i \setminus T_i$. First, we read the string $w = a_i^{f_i-s}$ which sends s to f_i in A_i and fixes all states in all A_j with $j \neq i$ to get states

$$\begin{aligned} (S'_1, S'_2, S'_3, \dots, S'_{i-1}, S' \cup (S_i \cdot w), S'_{i+1}, \dots, S'_k) \\ (T'_1, T'_2, T'_3, \dots, T'_{i-1}, T' \cup (T_i \cdot w), T'_{i+1}, \dots, T'_k) \end{aligned}$$

where $S', T' \subseteq [1, f_i - s]$ and $f_i \in (S_i \cdot w) \setminus (T_i \cdot w)$, that is, the i th components of the resulting states differ in the state f_i . If $S'_{i+1} \neq T'_{i+1}$, then we have the desired result. Otherwise, since $s_{i+1} \in S'_{i+1}$, both S'_{i+1} and T'_{i+1} are non-empty, which means that all S'_1, S'_2, \dots, S'_i and all T'_1, T'_2, \dots, T'_i are non-empty. Now, the string b sends all states of Q_j with $2 \leq j \leq k-1$, either to s_j or to f_j , and then $a_j b$ sends f_j to s_j and s_j to itself since $n_j \geq 3$. Thus after reading the string $b(a_2 b)(a_3 b) \cdots (a_{i-1} b)$ and if $T'_1 = \{f_1\}$, then also $(a_1 b)$, we get states

$$\begin{aligned} (\{q\}, \{s_2\}, \{s_3\}, \dots, \{s_{i-1}\}, S'' \cup \{f_i\}, S''_{i+1}, \dots, S''_k) \\ (\{s_1\}, \{s_2\}, \{s_3\}, \dots, \{s_{i-1}\}, \{s_i\}, T''_{i+1}, \dots, T''_k) \end{aligned}$$

where $q \in \{s_1, f_1\}$, $S'' \subseteq \{s_i\}$, and $S''_j, T''_j \subseteq \{s_j, f_j\}$ for $j = i+1, i+2, \dots, k-1$. There are two cases to consider.

Case 1: $1 \leq i \leq k-2$. Then $2 \leq i+1 \leq k-1$ and $n_{i+1} \geq 3$ which means that the string $a_{i+1}b$ sends both f_{i+1} and s_{i+1} to s_{i+1} . Thus after reading $a_{i+1}b$, we get states

$$\begin{aligned} &(\{q\}, \{s_2\}, \{s_3\}, \dots, \{s_{i-1}\}, S'' \cup \{f_i\}, \{s_{i+1}\}, S'''_{i+2}, \dots, S'''_k) \\ &(\{s_1\}, \{s_2\}, \{s_3\}, \dots, \{s_{i-1}\}, \{s_i\}, \{s_{i+1}\}, T'''_{i+2}, \dots, T'''_k). \end{aligned}$$

Finally, the string a_{i+1} , which performs the identity on Q_j with $j \neq i+1$ and the circular shift on Q_{i+1} , sends the resulting states to states which differ in s_{i+1} .

Case 2: $i = k-1$. Then the string b^{n_k} sends all states of Q_k to s_k , while it fixes s_j and f_j for $j = 1, 2, \dots, k-1$. Thus after reading the string b^{n_k} we get states $(\{q\}, \{s_2\}, \dots, \{s_{k-2}\}, S'' \cup \{f_{k-1}\}, \{s_k\})$ and $(\{s_1\}, \{s_2\}, \dots, \{s_{k-2}\}, \{s_{k-1}\}, \{s_k\})$. Now, in the same way as in Example 14 we show that either the string $a_{k-1}^{n_k}$ or the string $a_{k-1}^{n_k-1}ba_{k-1}^{n_k-1}$ sends the resulting states to two states which differ in s_k . \square

Since the number of valid states provides an upper bound on the state complexity of multiple concatenation, we get our main result.

Corollary 16. *The DFAs A_1, A_2, \dots, A_k shown in Figure 10 defined over a k -letter alphabet are witnesses for multiple concatenation of k languages.* \square

We conjecture that k symbols are necessary for describing witnesses for concatenation of k languages. The next observation shows that our conjecture holds for $k = 3$.

Theorem 17. *The ternary alphabet used to describe witnesses for the concatenation of three languages in Theorem 15 is optimal.*

Proof. Let $\Sigma = \{a, b\}$ and $n_i \geq 2$ for $i = 1, 2, 3$. Let us consider binary DFAs $A_i = (Q_i, \Sigma, \cdot, s_i, \{f_i\})$ where $Q_i = \{1, 2, \dots, n_i\}$, $s_i = 1$, $f_i \neq 1$ for $i = 1, 2, 3$; notice that to meet the upper bound for multiple concatenation, each A_1, A_2, \dots, A_{k-1} must have one final state, and it must be different from the initial state.

Construct the NFA N for $L(A_1)L(A_2)L(A_3)$ from DFAs A_1, A_2, A_3 as follows: for $i = 1, 2$, each state $q \in Q_i$ and each symbol $\sigma \in \{a, b\}$ such that $q\sigma = f_i$, add the transition (q, σ, s_{i+1}) ; the initial state of N is s_1 and its unique final state is f_3 . Our aim is to show that either some valid state is unreachable in the subset automaton $\mathcal{D}(N)$ or some valid states are equivalent to each other.

Notice that to reach the valid state $(s_1, Q_2, \{s_3\})$, we must have an input symbol that performs a permutation on Q_2 , and to reach the valid state $(s_1, \{s_2\}, Q_3)$, we must have an input symbol that performs a permutation on Q_3 .

If both input symbols perform a permutation on Q_3 , then the valid states $(s_1, \{s_2\}, Q_3)$ and $(s_1, \{2\}, Q_3)$ are equivalent since all strings are accepted from both of them.

If both input symbols perform a permutation on Q_2 , then the valid states $(s_1, Q_2, \{s_3\})$ and $(2, Q_2, \{s_3\})$ are equivalent since if a string w is accepted by N from the state s_1 in A_1 through a computation $s_1 \xrightarrow{w'} s_2 \xrightarrow{w''} f_3$ with $w = w'w''$, then it is accepted through a computation $w's_2 \xrightarrow{w'} s_2 \xrightarrow{w''} f_3$ where $w's_2 \in Q_2$, so it is accepted from $(2, Q_2, \{s_3\})$; and vice versa.

Hence to meet the upper bound, we must have one permutation and one non-permutation input symbol in both A_2, A_3 .

Next, while reaching the valid state $(s_1, Q_2 \setminus \{f_2\}, \emptyset)$, we cannot visit state f_2 . This means that there must be an input that maps $Q_2 \setminus \{f_2\}$ onto $Q_2 \setminus \{f_2\}$. Without loss of generality, let this input be a . Since f_2 must be reachable in A_2 , there must exist a state p in $Q_2 \setminus \{f_2\}$ with $pb = f_2$. Moreover, $f_2b \neq f_2$ because otherwise either f_2 would have loops on both symbols, or both a and b would be non-permutation symbols in A_2 . We have two cases:

(1) Let b be a non-permutation symbol in A_2 . Then a is a permutation on Q_2 , so $f_2a = f_2$. This situation is depicted in Fig. 11. Moreover, there is a state in $Q_2 \setminus \{f_2\}$ with no in-transition on b . Therefore the valid state $(s_1, Q_2 \setminus \{f_2\}, Q_3)$ must be reached from some valid state on a , and consequently a is a permutation on Q_3 . Next, since $f_2b \neq f_2$, the valid state $(s_1, \{f_2b\}, Q_3)$ must be reached from a valid state $(j, \{f_2\} \cup S, Q_3)$ on b since to get Q_3 in the third component, we must visit f_2 , and only reading b eliminates the state f_2 . It follows that b is a permutation on Q_3 . Hence both a and b perform permutations on Q_3 , thus resulting in a contradiction.

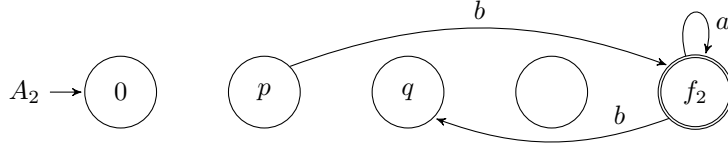


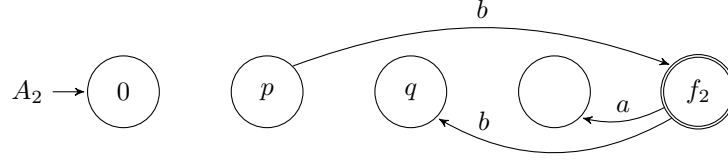
Figure 11: Case 1: a maps $Q_2 \setminus \{f_2\}$ onto $Q_2 \setminus \{f_2\}$ and b is not a permutation on Q_2 .

(2) Let b be a permutation symbol in A_2 . Then a is not a permutation on Q_2 , so $f_2a \neq f_2$, and therefore $f_2 \notin Q_2a$, so each state containing f_2 in its second component must be reached by b . This situation is illustrated in Fig. 12. It follows that every valid state $(j, Q_2, \{s_3\})$ must be reached on b , so b is a permutation on Q_1 ,

Next, the valid state $(s_1, \{f_2\}, Q_3)$ must be reached on b as well. Therefore each state in $Q_3 \setminus \{s_3\}$ has an in-transition on b . Moreover, the state $(f_1b, Q_2, \{s_3\})$ must be reached by b from a valid state $(f_1, Q_2, \{s_3\} \cup T)$; recall that b is a permutation on Q_1 . This means that $s_3b = s_3$. Hence b is a permutation on Q_3 . Let $r \in Q_2 \setminus \{s_2b, f_2\}$. Then the valid state $(f_1b, \{r\}, Q_3)$ cannot be reached on b because otherwise it would be reached from $(f_1, \{s_2\} \cup S, T)$ and would contain s_2b in its second component. It follows that a is a permutation on Q_3 . Thus both a and b perform a permutation in A_2 , which is a contradiction. \square

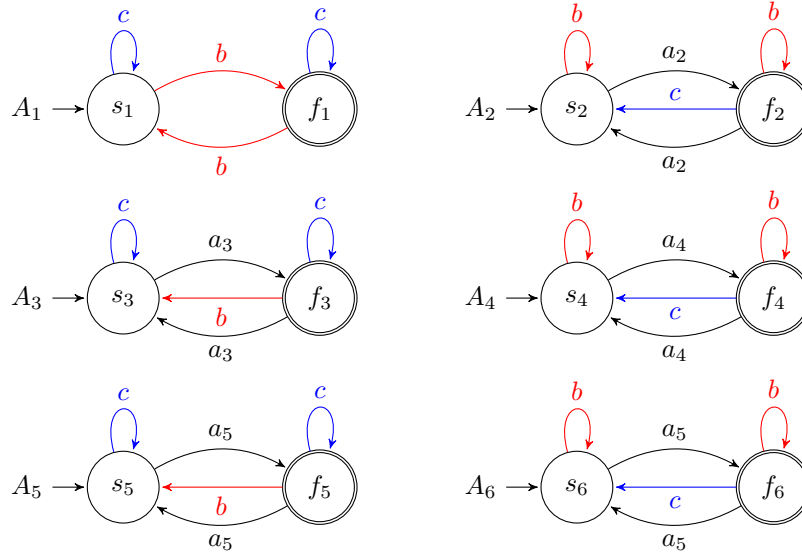
Notice that all our k -letter witness DFAs from Theorem 15, except for the first and last one, are assumed to have at least three states. However, our witnesses over a $(k+1)$ -letter alphabet from Theorem 13 cover also the cases when some of given DFAs have two states. Although, we are not able to cover such cases by using just k letters, we can do it providing that all automata have two states. We only give the main ideas here for this case.

Let $\Sigma = \{b, c, a_2, a_3, \dots, a_{k-1}\}$ be a k -letter alphabet. For $i = 1, 2, \dots, k$, let $A_i = (Q_i, \Sigma, s_i, \cdot, f_i)$ be a two-state DFA with $Q_i = \{1, 2\}$, $s_i = 1$, $f_i = 2$,

Figure 12: Case 2: a maps $Q_2 \setminus \{f_2\}$ onto $Q_2 \setminus \{f_2\}$ and b is a permutation on Q_2 .

and the transitions defined as follows (see Figure 13 for an illustration):

- a_i with $i = 2, 3, \dots, k-2$ performs the cycle on Q_i and the identity on Q_j with $j \neq i$;
- a_{k-1} performs the cycle on Q_{k-1} and Q_k , and the identity on Q_1, Q_2, \dots, Q_{k-2} ;
- b performs the cycle on Q_1 , the identity on Q_i if i is even, and the contraction ($f_i \rightarrow s_i$) on Q_i if $i \geq 3$ is odd;
- c performs the identity on Q_i if i is odd, and the contraction ($f_i \rightarrow s_i$) otherwise.

Figure 13: Two-state DFAs; $k = 6$. In each DFA, the remaining symbols perform identities.

Construct an NFA N for $L(A_1)L(A_2)\cdots L(A_k)$ from the DFAs A_1, A_2, \dots, A_k as follows: for each $i = 1, 2, \dots, k-1$, each $q \in Q_i$ and $\sigma \in \Sigma$ such that $q \cdot \sigma = f_i$ in A_i , add the transition (q, σ, s_{i+1}) ; the initial state of N is s_1 and its final state is f_k .

We prove reachability and distinguishability of states of the subset automaton $\mathcal{D}(N)$ in a similar way as before, but we have take into account that to reach a state $p = (f_1, T_2, T_3, \dots, T_k)$ from a state $q = (s_1, S_2, S_3, \dots, S_k)$, the symbol b has

to be read. However, although b sends s_1 to f_1 , it also sends each non-empty subset S_i with $i \geq 3$ and i odd to $\{s_i\}$. Then, we have to carefully return $\{s_i\}$ back to S_i .

6. Binary and Ternary Languages

In this section, we examine the state complexity of multiple concatenation on binary and ternary languages. Our aim is to show that in the binary case, the resulting complexity is still exponential in n_2, n_3, \dots, n_k , and in the ternary case, it is the same as in the general case, up to a multiplicative constant depending on k . Let us start with the following example.

Example 18. Let $n \geq 3$ and N be the NFA shown in Figure 14 that recognizes the language of strings over $\{a, b\}$ which have an a in the $(n-1)$ st position from the end.

Let us show that each subset $S \subseteq [1, n]$ with $1 \in S$ is reachable in the subset automaton $\mathcal{D}(N)$. The proof is by induction on $|S|$. The basis, with $|S| = 1$, holds true since $\{1\}$ is the initial state. Let $|S| \geq 2$ and $1 \in S$. Let $m = \min(S \setminus \{1\})$. Set $S' = ab^{m-2}(S \setminus \{1, m\})$. Then $S' \subseteq [2, n - m + 1]$ and $|S'| = |S| - 2$. We have $\{1\} \cup S' \xrightarrow{a} \{1, 2\} \cup b^{s-2}(S \setminus \{1, s\}) \xrightarrow{b^{s-2}} \{1, s\} \cup (S \setminus \{1, s\}) = S$, where the leftmost set of size $|S| - 1$ is reachable by induction. ■

We now use the result from the above example to get a lower bound on the state complexity of multiple concatenation on binary languages. The idea is to describe binary DFAs in such a way that the NFA for their concatenation would accept, except for a finite set, the set of strings having an a in an appropriate position from the end.

Theorem 19. Let $k \geq 3$, $n_1 \geq 3$, $n_2 \geq 4$, and $n_i \geq 3$ for $i = 3, 4, \dots, k$. Let A_1, A_2, \dots, A_k be the binary DFAs shown in Figure 15. Then every DFA for the language $L(A_1)L(A_2) \cdots L(A_k)$ has at least $n_1 - 1 + (1/2^{2k-2})2^{n_2+n_3+\dots+n_k}$ states.

Proof. Construct an NFA for $L(A_1)L(A_2) \cdots L(A_k)$ from the DFAs A_1, A_2, \dots, A_k by adding the transitions (f_1-1, b, s_2) , (f_1, a, s_2) , (f_1, b, s_2) , and (f_i-1, σ, s_{i+1}) for $i = 2, 3, \dots, k-1$ and $\sigma \in \{a, b\}$, by making states f_1, f_2, \dots, f_{k-1} non-final, and states s_2, s_3, \dots, s_k non-initial. In this NFA, the states f_i and f_i+1 with $2 \leq i \leq k-1$, as well as the state f_k+1 are dead, so we can omit them. Let N be the resulting NFA; see Figure 16 for an illustration.

In the subset automaton $\mathcal{D}(N)$, each state $(j, \emptyset, \emptyset, \dots, \emptyset)$ with $1 \leq j \leq f_1 - 1$ is reached from the initial state $(s_1, \emptyset, \emptyset, \dots, \emptyset)$ by b^{j-1} , and $(f_1, \{s_2\}, \emptyset, \emptyset, \dots, \emptyset)$ is reached from $(f_1-1, \emptyset, \emptyset, \dots, \emptyset)$ by b . Starting with the state f_1 , the NFA N accepts

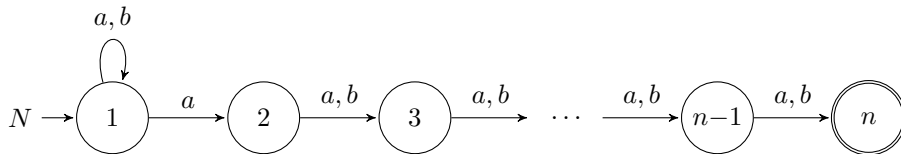


Figure 14: A binary NFA N such that every set $\{1\} \cup S$ is reachable in $\mathcal{D}(N)$.

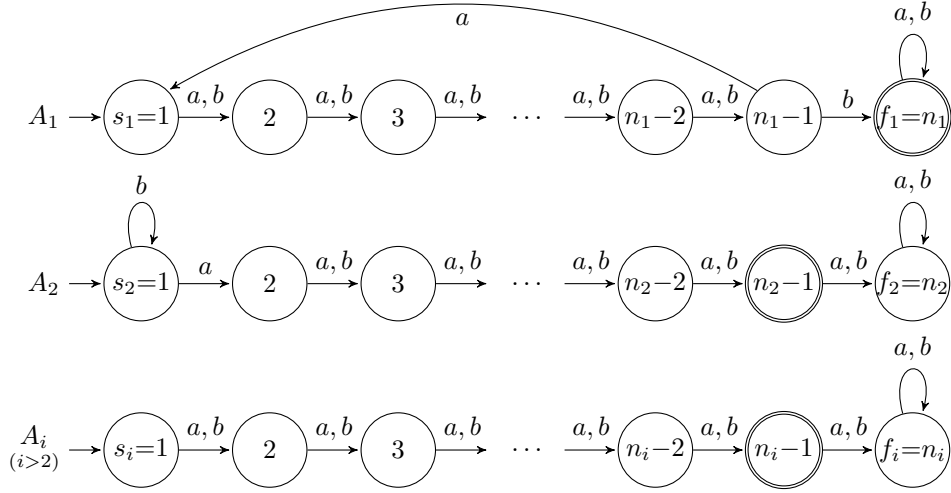


Figure 15: Binary DFAs A_1, A_2 , and A_i for $i = 3, 4, \dots, k$ meeting the lower bound $n_1 - 1 + (1/2^{2k-1})2^{n_2+n_3+\dots+n_k}$ for multiple concatenation.

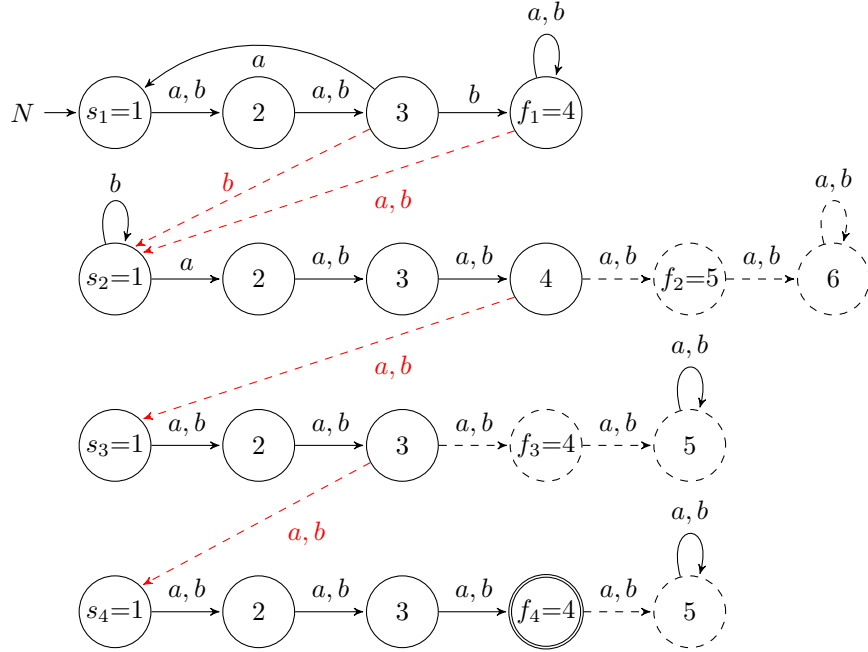


Figure 16: A binary NFA for $L(A_1)L(A_2)L(A_3)L(A_4)$ where $n_1 = 4, n_2 = 6, n_3 = n_4 = 5$.

all strings having an a in position $n_2 - 2 + n_3 - 2 + \dots + n_{k-1} - 2 + n_k - 1$ from the end. As shown in Example 18, every state $(f_1\{s_2\} \cup S_2, S_3, \dots, S_k)$ with $S_2 \subseteq \{2, 3, \dots, n_2 - 2\}$, $S_i \subseteq \{1, 2, \dots, n_i - 2\}$ for $i = 3, 4, \dots, k - 1$, and $S_k \subseteq \{1, 2, \dots, n_k - 1\}$ is reachable. This gives $n_1 - 1 + 2^{n_2 - 3 + n_3 - 2 + \dots + n_{k-1} - 2 + n_k - 1} = n_1 - 1 + (1/2^{2k-2})2^{n_2 + n_3 + \dots + n_k}$ reachable states.

Moreover, each singleton set is co-reachable in N via a string in a^* , except for $\{q\}$ where q is a non-final state of A_1 . By Lemma 1, the reachable states $(i, S_2, S_3, \dots, S_k)$ and $(j, T_2, T_3, \dots, T_k)$ are distinguishable if they differ in a state of A_i with $i \geq 2$ or in f_1 . Next, the states $(i, S_2, S_3, \dots, S_k)$ and $(j, S_2, S_3, \dots, S_k)$ with $1 \leq i < j < f_1$ are sent to states that differ in f_1 by b^{f_1-j} . \square

Our next result shows that a trivial upper bound $n_1 2^{n_2 + n_3 + \dots + n_k}$ can be met, up to a multiplicative constant depending on k , by the concatenation of k ternary languages. Thus, this trivial upper bound is asymptotically tight in the ternary case.

Theorem 20. *Let $k \geq 2$, $n_1 \geq 3$, $n_2 \geq 4$, and $n_i \geq 3$ for $i = 3, 4, \dots, k$. There exist ternary DFAs A_1, A_2, \dots, A_k such that every DFA recognizing the concatenation $L(A_1)L(A_2) \dots L(A_k)$ has at least $(1/2^{2k-2})n_1 2^{n_2 + n_3 + \dots + n_k}$ states.*

Proof. Let us add the transitions on symbol c to the binary automata shown in Figure 15 as follows: $c: (1, 2, \dots, n_1)$ in A_1 , $c: (f_i \rightarrow f_i + 1)$ in A_i with $2 \leq i \leq k - 1$, and $c: (1)$ in A_k . Construct the NFA N for $L(A_1)L(A_2) \dots L(A_k)$ with omitted dead states as in the binary case; see Figure 17 for an illustration. As shown in the proof of Theorem 19, the subset automaton $\mathcal{D}(N)$ has $(1/2^{2k-2})2^{n_2 + n_3 + \dots + n_k}$ reachable states of the form $(f_1, S_2, S_3, \dots, S_k)$. Each such state is sent to the state $(j, S_2, S_3, \dots, S_k)$ with $1 \leq j \leq f_1 - 1$ by the string c^j . Moreover, in the NFA N , each singleton set is co-reachable via a string in a^*c^* . By Corollary 2, all states of $\mathcal{D}(N)$ are pairwise distinguishable. This gives the desired lower bound. \square

7. Unary Languages

The upper bound on the state complexity of concatenation of two unary languages is $n_1 n_2$, and this upper bound can be met by cyclic unary languages if $\gcd(n_1, n_2) = 1$ as shown in [8, Theorems 5.4 and 5.5]. This gives a trivial upper bound $n_1 n_2 \dots n_k$ for concatenation of k unary languages. Here we show that a tight upper bound for concatenation of k cyclic unary languages is much smaller. Then we continue our study by investigating the concatenation of languages of the form $a^{\mu_i} Y_i$ where Y_i is a λ_i -cyclic. In both cases, we provide tight upper bounds. Finally, we consider the case, when automata may have final states in their tails.

Recall that the state set of a unary automaton of size (λ, μ) consists of a tail $q_0, q_1, \dots, q_{\mu-1}$ and a cycle $p_0, p_1, \dots, p_{\lambda-1}$ (with $p_0 = q_0$ if $\mu = 0$), and its transitions are $q_0 \rightarrow q_1 \rightarrow \dots \rightarrow q_{\mu-1} \rightarrow p_0 \rightarrow p_1 \rightarrow \dots \rightarrow p_{\lambda-1} \rightarrow p_0$; cf. [6].

Let n_1, n_2, \dots, n_k be positive integers with $\gcd(n_1, n_2, \dots, n_k) = 1$. Then $g(n_1, n_2, \dots, n_k)$ denotes the Frobenius number, that is, the largest integer that cannot be expressed as $x_1 n_1 + x_2 n_2 + \dots + x_k n_k$ for some non-negative integers x_1, x_2, \dots, x_k . Let us start with the following observation.

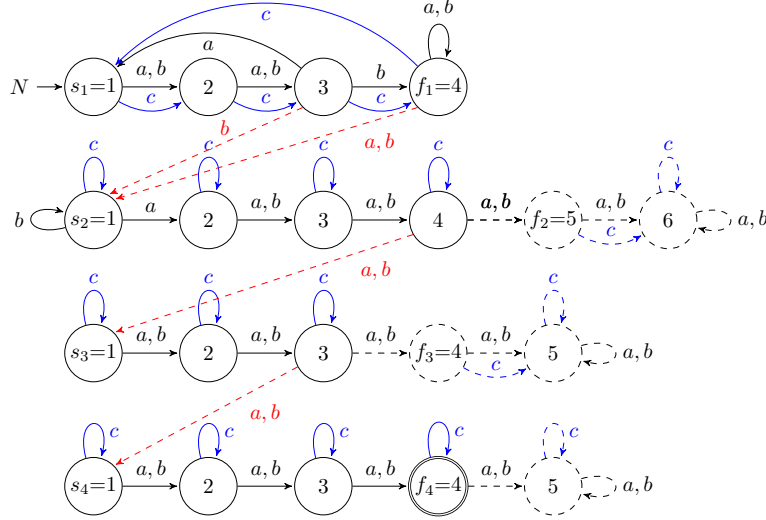


Figure 17: A ternary NFA for $L(A_1)L(A_2)L(A_3)L(A_4)$ where $n_1 = 4$, $n_2 = 6$, $n_3 = n_4 = 5$.

Lemma 21. *Let n_1, n_2, \dots, n_k be positive integers with $\gcd(n_1, n_2, \dots, n_k) = d$. Then each number of the form $x_1n_1 + x_2n_2 + \dots + x_kn_k$, with $x_1, x_2, \dots, x_k \geq 0$, is a multiple of d . Furthermore, the largest multiple of d that cannot be represented as $x_1n_1 + x_2n_2 + \dots + x_kn_k$, with $x_1, x_2, \dots, x_k \geq 0$, is $d \cdot g(\frac{n_1}{d}, \frac{n_2}{d}, \dots, \frac{n_k}{d})$.*

Proof. The first claim follows from the fact that each n_i is a multiple of d . Since $\gcd(\frac{n_1}{d}, \frac{n_2}{d}, \dots, \frac{n_k}{d}) = 1$, the largest integer that cannot be represented as $x_1\frac{n_1}{d} + x_2\frac{n_2}{d} + \dots + x_k\frac{n_k}{d}$, with $x_1, x_2, \dots, x_k \geq 0$, is $g(\frac{n_1}{d}, \frac{n_2}{d}, \dots, \frac{n_k}{d})$. Multiplying by d , we get the second claim. \square

Let $f(n_1, n_2, \dots, n_k) = g(n_1, n_2, \dots, n_k) + n_1 + n_2 + \dots + n_k$ be the modified Frobenius number, that is, the largest integer which is not representable by positive integer linear combinations. Using this notation, we have the following result.

Theorem 22. *Let A_1, A_2, \dots, A_k be unary cyclic automata with n_1, n_2, \dots, n_k states, respectively. Let $d = \gcd(n_1, n_2, \dots, n_k)$. Then $L(A_1)L(A_2) \dots L(A_k)$ is recognized by a DFA of size (λ, μ) , where $\lambda = d$ and $\mu = d \cdot f(\frac{n_1}{d}, \frac{n_2}{d}, \dots, \frac{n_k}{d}) - k + 1$, and this upper bound is tight.*

Proof. Denote $L_i = L(A_i)$ and $L = L_1L_2 \dots L_k$. We show that L is recognized by a unary DFA of size (λ, μ) . By [6, Theorem 2], it is enough to show that for every $m \geq d \cdot f(\frac{n_1}{d}, \frac{n_2}{d}, \dots, \frac{n_k}{d}) - k + 1$, we have $a^m \in L$ if and only if $a^{m+d} \in L$.

We can write each language L_i as $L_i = Z_i(a^{n_i})^*$ where $Z_i = L_i \cap \{a^x \mid 0 \leq x < n_i\}$; cf. [6, Proof of Theorem 8]. Let $m \geq d \cdot f(\frac{n_1}{d}, \frac{n_2}{d}, \dots, \frac{n_k}{d}) - k + 1$.

If $a^m \in L$, then $m = z_1 + x_1n_1 + z_2 + x_2n_2 + \dots + z_k + x_kn_k$ where $a^{z_i} \in Z_i$

and $x_i \geq 0$. Since $m \geq d \cdot f(\frac{n_1}{d}, \frac{n_2}{d}, \dots, \frac{n_k}{d}) - k + 1$, we get

$$\begin{aligned} x_1 n_1 + x_2 n_2 + \dots + x_k n_k &\geq d \cdot f(\frac{n_1}{d}, \frac{n_2}{d}, \dots, \frac{n_k}{d}) - k + 1 - z_1 - z_2 - \dots - z_k \geq \\ &d \cdot f(\frac{n_1}{d}, \frac{n_2}{d}, \dots, \frac{n_k}{d}) - k + 1 - (n_1 - 1) - (n_2 - 1) - \dots - (n_k - 1) = \\ &d \cdot g(\frac{n_1}{d}, \frac{n_2}{d}, \dots, \frac{n_k}{d}) + 1. \end{aligned}$$

Since $x_1 n_1 + x_2 n_2 + \dots + x_k n_k$ is a multiple of d , it follows from Lemma 21 that $x_1 n_1 + x_2 n_2 + \dots + x_k n_k + d = x'_1 n_1 + x'_2 n_2 + \dots + x'_k n_k$ for some $x'_1, x'_2, \dots, x'_k \geq 0$. Therefore

$$m + d = z_1 + x'_1 n_1 + z_2 + x'_2 n_2 + \dots + z_k + x'_k n_k,$$

so $a^{m+d} \in L$.

Conversely, if $a^{m+d} \in L$, then $m + d = z_1 + x_1 n_1 + z_2 + x_2 n_2 + \dots + z_k + x_k n_k$ where $a^{z_i} \in Z_i$ and $x_i \geq 0$. Since $m \geq d \cdot f(\frac{n_1}{d}, \frac{n_2}{d}, \dots, \frac{n_k}{d}) - k + 1$, similarly as in the previous paragraph, we get

$$x_1 n_1 + x_2 n_2 + \dots + x_k n_k - d \geq d \cdot g(\frac{n_1}{d}, \frac{n_2}{d}, \dots, \frac{n_k}{d}) + 1,$$

and therefore $x_1 n_1 + x_2 n_2 + \dots + x_k n_k - d = x'_1 n_1 + x'_2 n_2 + \dots + x'_k n_k$ for some $x'_1, x'_2, \dots, x'_k \geq 0$. Thus $m = z_1 + x'_1 n_1 + z_2 + x'_2 n_2 + \dots + z_k + x'_k n_k$ and $a^m \in L$.

To get tightness, consider unary cyclic languages $L_i = a^{n_i-1}(a^{n_i})^*$ recognized by unary cyclic n_i -state automata. Let $L = L_1 L_2 \dots L_k$. As shown above, the language L is recognized by a unary DFA A with a tail of length $d \cdot f(\frac{n_1}{d}, \dots, \frac{n_k}{d}) - k + 1$ and a cycle of size d . Next, we have $a^m \in L$ if and only if

$$m = (n_1 - 1) + (n_2 - 1) + \dots + (n_k - 1) + x_1 n_1 + x_2 n_2 + \dots + x_k n_k$$

for some $x_1, x_2, \dots, x_k \geq 0$. Since $x_1 n_1 + x_2 n_2 + \dots + x_k n_k$ is a multiple of d , the cycle of size d has exactly one final state, and therefore it is minimal. Furthermore, a string $a^{d \cdot f(\frac{n_1}{d}, \frac{n_2}{d}, \dots, \frac{n_k}{d}) - k + \ell d}$ is in L if and only if

$$d \cdot f(\frac{n_1}{d}, \frac{n_2}{d}, \dots, \frac{n_k}{d}) - k + \ell d = (n_1 - 1) + (n_2 - 1) + \dots + (n_k - 1) + x_1 n_1 + x_2 n_2 + \dots + x_k n_k$$

for some $x_1, x_2, \dots, x_k \geq 0$, which holds if and only if

$$d \cdot g(\frac{n_1}{d}, \frac{n_2}{d}, \dots, \frac{n_k}{d}) + \ell d = x_1 n_1 + x_2 n_2 + \dots + x_k n_k.$$

By Lemma 21, it follows that $a^{d \cdot f(\frac{n_1}{d}, \frac{n_2}{d}, \dots, \frac{n_k}{d}) - k} \notin L$, while $a^{d \cdot f(\frac{n_1}{d}, \frac{n_2}{d}, \dots, \frac{n_k}{d}) - k + d} \in L$. Hence A is minimal. \square

By [3, Proposition 2.2], if $n_1 \leq n_2 \leq \dots \leq n_k$, then $g(n_1, n_2, \dots, n_k) \leq n_1 n_k$. This gives an upper bound $n_1 n_k / d + n_1 + \dots + n_k - k + 1 + d$ for concatenation of k cyclic languages where $n_1 \leq n_2 \leq \dots \leq n_k$ and $d = \gcd(n_1, n_2, \dots, n_k)$. The result of the previous theorem can be generalized as follows.

Corollary 23. For $i = 1, 2, \dots, L_k$, let $L_i = a^{\mu_i} Y_i$ where Y_i is λ_i -cyclic be a language recognized by a DFA of size (λ_i, μ_i) . Let $d = \gcd(\lambda_1, \lambda_2, \dots, \lambda_k)$. Then the language $L_1 L_2 \dots L_k$ is recognized by a DFA of size (λ, μ) where $\lambda = d$ and $\mu = \mu_1 + \mu_2 + \dots + \mu_k + d \cdot f(\frac{\lambda_1}{d}, \frac{\lambda_2}{d}, \dots, \frac{\lambda_k}{d}) - k + 1$, and this upper bound is tight.

Proof. The language $L_1 L_2 \dots L_k$ is a concatenation of the singleton language $a^{\mu_1 + \mu_2 + \dots + \mu_k}$ recognized by a DFA of size $(1, \mu_1 + \mu_2 + \dots + \mu_k + 1)$ and the concatenation of cyclic languages $Y_1 Y_2 \dots Y_k$. Now the result follows from the previous theorem since we can simply merge the final state of the automaton for the singleton language with the initial state of the DFA for $Y_1 Y_2 \dots Y_k$; cf. [6, Theorem 6]. The upper bound is met by languages $L_i = a^{\mu_i + \lambda_i - 1} (a^{\lambda_i})^*$. \square

In the case of concatenation of two languages, the length of the resulting cycle may be equal to the least common multiple of the lengths of cycles in given automata providing that they have final states in their tails [6, Theorems 10 and 11]. The next example shows that in some cases this is the optimal way how to get the maximum complexity of concatenation of languages recognized by m -state and n -state unary DFAs, respectively.

Example 24. Given an m -state and n -state unary DFA, their concatenation requires mn states if $\gcd(m, n) = 1$. If $\gcd(m, n) > 1$, then we may try to take DFAs with smaller cycles of sizes $m-i$ and $n-j$, and inspect the complexity of concatenation of languages recognized by automata of sizes $(m-i, i)$ and $(n-j, j)$.

As shown in [6, Theorem 11] the minimal DFA for concatenation of the languages $\{\varepsilon\} \cup a^{m-1} (a^{m-2})^*$ and $\{\varepsilon\} \cup a^{n-1} (a^{n-2})^*$, that are recognized by automata of sizes $(m-2, 2)$ and $(n-2, 2)$, with the set of final states $\{0, m-1\}$ and $\{0, n-1\}$, respectively, has $2\text{lcm}(m-2, n-2) + 3$ states. By our computations, the smallest m and n , for which such automata provide the maximum complexity among all automata of sizes $(m-i, i)$ and $(n-j, j)$, are $m = 137\,712$ and $n = 127\,206$.

Nevertheless, it looks like sometimes it could be helpful to decrease the lengths of cycles not by two, but just by one, and setting the final state sets to $\{0, m-2\}$ and $\{0, n-2\}$, respectively; our aim is to have a state in both tails, and then, to get minimal DFAs, the states $m-1$ and $n-1$ have to be non-final. Then, similarly as in the proof of [6, Theorem 11] we show that the minimal DFA recognizing the concatenation of these two languages has $2\text{lcm}(m-1, n-1) - 1$ states provided that $\gcd(m-1, n-1) > 1$ and neither $m-1$ nor $n-1$ is a multiple of the other.

Our next goal is to find m and n such that the maximum of complexities of concatenation of languages recognized by all automata of sizes $(m-i, i)$ and $(n-j, j)$ is achieved if $i = j = 1$ and $\gcd(m-1, n-1) = 2$ by the above mentioned languages. In such a case, we have $2\text{lcm}(m-1, n-1) - 1 = (m-1)(n-1) - 1$.

By [6, Theorems 10 and 12], the complexity of concatenation of languages recognized by automata of sizes $(m-i, i)$ and $(n-j, j)$ is at most $(m-i)(n-j) + i + j$ if $\gcd(m-i, n-j) = 1$, and at most $2\text{lcm}(m-i, n-j) + i + j - 1$ if $\gcd(m-i, n-j) > 1$. In both cases, the resulting complexity is at most $(m-i)(n-j) + i + j$. Denote this

number by $c_{i,j} = (m-i)(n-j) + i + j$. The reader may verify that

$$\begin{aligned} c_{i,j} &< (m-1)(n-1) - 1 \text{ for all } i, j \geq 1 \text{ and } (i, j) \neq (1, 1), \\ c_{0,j} &< (m-1)(n-1) - 1 \text{ if } j \geq 2 \text{ and } n+2 < m, \\ c_{i,0} &< (m-1)(n-1) - 1 \text{ if } i \geq 3 \text{ and } m < 2n-3. \end{aligned}$$

It follows that the complexity $(m-1)(n-1) - 1$ could possibly be exceeded only by automata of sizes $(m-i, i)$ and $(n-j, j)$ where $(i, j) \in \{(0, 0), (0, 1), (1, 0), (2, 0)\}$. Assume that in all of these cases, we have $\gcd(m-i, n-j) \geq 3$. Then, providing that $m, n \geq 8$, the complexity of the corresponding concatenations in these four cases is at most

$$2 \operatorname{lcm}(m-i, n-j) + i + j - 1 < \frac{2}{3}(m-i)(n-j) + i + j \leq \frac{2}{3}mn + 3 < (m-1)(n-1) - 1.$$

Now, let $m = 471$ and $n = 315$. Then $\gcd(m-1, n-1) = 2$ and $n+2 < m < 2n-3$. Moreover, we have $\gcd(471, 315) = 3$, $\gcd(471, 314) = 157$, $\gcd(470, 315) = 5$, and $\gcd(469, 315) = 7$. This means that the maximum complexity of concatenation of a 471-state and 315-state unary DFA is achieved by automata of sizes $(470, 1)$ and $(314, 1)$ recognizing languages $\{\varepsilon\} \cup a^{469}(a^{470})^*$ and $\{\varepsilon\} \cup a^{313}(a^{314})^*$, that is, by automata that have a final state in their tails. ■

Motivated by our previous examples, we finally consider the state complexity of the concatenation of k languages recognized by unary automata that have final states in their tails. While in our previous two theorems, the length of the resulting cycle was equal to the greatest common divisor of the lengths of cycles in the given automata, here, similarly to the case of concatenation of two languages (cf. [6, Theorems 10, 11]), it may be equal to their least common multiple. We cannot obtain a tight upper bound here, nevertheless, we provide an example that meets our upper bound.

Theorem 25. *For $i = 1, 2, \dots, k$, let A_i be a unary DFA of size (λ_i, μ_i) . For a non-empty set $I = \{i_1, i_2, \dots, i_\ell\} \subseteq \{1, 2, \dots, k\}$, let*

$$\begin{aligned} d_I &= \gcd(\lambda_{i_1}, \lambda_{i_2}, \dots, \lambda_{i_\ell}), \\ f(I) &= f\left(\frac{\lambda_{i_1}}{d_I}, \frac{\lambda_{i_2}}{d_I}, \dots, \frac{\lambda_{i_\ell}}{d_I}\right), \end{aligned}$$

and set $d_\emptyset = 1$ and $f(\emptyset) = 0$. Then the language $L(A_1)L(A_2) \cdots L(A_k)$ is recognized by a DFA of size (λ, μ) where

$$\begin{aligned} \lambda &= \operatorname{lcm}(\lambda_1, \lambda_2, \dots, \lambda_k) \\ \mu &= \max\{\mu_1 + \mu_2 + \dots + \mu_k - k + 1 + d_I \cdot f(I) \mid I \subseteq \{1, 2, \dots, k\}\}. \end{aligned}$$

Proof. Let $L_i = L(A_i)$ and $L = L(A_1)L(A_2) \cdots L(A_k)$. We have $L_i = X_i \cup a^{\mu_i}Y_i$ where $X_i = L(A_i) \cap \{a^x \mid 0 \leq x < \mu_i\}$ and $Y_i = \{a^x \mid a^{\mu_i+x} \in L(A_i)\}$. Then

$$L = \bigcup_{I \subseteq \{1, 2, \dots, k\}} \prod_{j \notin I} X_j \prod_{i \in I} a^{\mu_i} Y_i.$$

For each I , the language $\prod_{j \notin I} X_j$ is a finite language recognized by a DFA of size $(1, 1 + \sum_{j \notin I} (\mu_j - 1))$, and by Corollary 23, the language $\prod_{i \in I} a^{\mu_i} Y_i$ is recognized by a DFA of size $(d_I, 1 + d_I \cdot F(I) + \sum_{i \in I} (\mu_i - 1))$.

The concatenation of these two languages is recognized by a DFA of size $(d_I, \mu_1 + \mu_2 + \dots + \mu_k - k + 1 + d_I \cdot f(I))$; cf. [6, Theorem 6]. Then, the union of these concatenations is recognized by a DFA of size (λ, μ) by [6, Theorem 4]. \square

Example 26. Consider unary DFAs A_1, A_2, A_3 of sizes $(12, 2)$, $(20, 2)$, and $(30, 2)$, with $F_1 = \{0, 13\}$, $F_2 = \{0, 21\}$, and $F_3 = \{0, 31\}$.

We have $\text{lcm}(12, 20, 30) = 60$, $4 \cdot f(3, 5) = 6 \cdot f(2, 5) = 10 \cdot f(2, 3) = 60$, and $2 \cdot f(6, 10, 15) = 2 \cdot 2 \cdot f(3, 5, 15) = 2 \cdot 2 \cdot 5 \cdot f(3, 1, 3) = 2 \cdot 2 \cdot 5 \cdot 3 \cdot f(1, 1, 1) = 2 \cdot 2 \cdot 5 \cdot 3 \cdot 2 = 120$. The size of the minimal automaton recognizing the language $L(A_1)L(A_2)L(A_3)$ is $(60, 124)$ where $124 = 2 + 2 + 2 - 3 + 1 + \max\{60, 120\}$. \blacksquare

The above example shows that our upper bound given by Theorem 25 is met by unary automata of sizes $(12, 2)$, $(20, 2)$, $(30, 2)$. The tightness of this upper bound in a general case remains open.

8. Conclusions

We examined in detail the state complexity of the multiple concatenation of k languages. First, we described witness DFAs A_1, A_2, \dots, A_k over the $(k + 1)$ -letter alphabet $\{b, a_1, a_2, \dots, a_k\}$, in which each a_i performs the circular shift in A_i and the identity in the other automata, while b performs a contraction. Using symbols a_1, a_2, \dots, a_k , we proved the reachability of all valid states in the subset automaton for the concatenation by carefully setting the i th component without changing the already set $(i + 1)$ th component. The transitions on b guaranteed the co-reachability of all singleton sets in the NFA for concatenation, and therefore we obtained the proof of distinguishability of all states in the corresponding subset automaton for free. However, to get co-reachability of singletons, our witness automata were required to have at least three states. Nevertheless, we described witness automata over a $(k + 1)$ -letter alphabet also in the case where some of them have only two states.

Then we provided special binary witnesses for the concatenation of two languages. Using our results concerning witnesses over a $(k + 1)$ -letter alphabet, as well as the results for the special binary automata, we described witnesses for the concatenation of k languages over a k -letter alphabet. This solves an open problem stated in [1]. For $k = 3$, we proved that the ternary alphabet is optimal in the sense that the upper bound for the concatenation of three languages cannot be met by any binary languages. This provides a partial answer to the second open problem from [1].

We also considered multiple concatenation on binary and ternary languages, and obtained lower bounds $n_1 - 1 + (1/2^{2k-2})2^{n_2+n_3+\dots+n_k}$ and $(1/2^{2k-2})n_12^{n_2+n_3+\dots+n_k}$, respectively. This shows that the state complexity of multiple concatenation remains exponential in n_2, n_3, \dots, n_k in the binary case, and that a trivial upper bound can be met, up to a multiplicative constant depending on k , by ternary languages.

Finally, we investigated multiple concatenation on unary languages. We obtained a tight upper bound for cyclic languages, and we showed that for $k \geq 3$, it is much smaller than a trivial upper bound $n_1 n_2 \cdots n_k$, which is met by cyclic unary languages if $k = 2$ and $\gcd(n_1, n_2) = 1$ [8, Theorem 5.4]. We also provided a tight upper bound for languages recognized by automata that do not have final states in their tails.

Some problems remain open. First, our k -letter witnesses require $n_i \geq 3$ for $i = 2, 3, \dots, k-1$, while the $(2k-1)$ -letter witnesses in [4, Theorem 5] work with $n_i \geq 2$. Is it possible to define k -letter witnesses also in such a case? We can do this using $k+1$ letters, or with k letters if *all* automata have two states.

We proved the optimality of a ternary alphabet for the concatenation of three languages. However, we cannot see any generalization of the proof. Is a k -letter alphabet for describing witnesses for the concatenation of k languages optimal?

Next, we provided upper bounds in the case where exactly one automaton has one state, and using a binary alphabet we proved that they are tight if $k = 2$. What is the state complexity of multiple concatenation if some languages may be equal to Σ^* ?

Finally, in the unary case, we obtained an upper bound for multiple concatenation of languages recognized by unary automata that may have final states in their tails. The tightness of this upper bound remains open.

References

- [1] P. CARON, J. LUQUE, B. PATROU, State complexity of multiple concatenations. *Fund. Inform.* **160** (2018) 3, 255–279.
<https://doi.org/10.3233/FI-2018-1683>
- [2] Z. ÉSIK, Y. GAO, G. LIU, S. YU, Estimation of state complexity of combined operations. *Theoret. Comput. Sci.* **410** (2009) 35, 3272–3280.
<https://doi.org/10.1016/j.tcs.2009.03.026>
- [3] J. GALLIER, The Frobenius coin problem upper bounds on the Frobenius number. Available online at: <https://www.cis.upenn.edu/~cis5110/Frobenius-number.pdf>.
- [4] Y. GAO, S. YU, State complexity approximation. In: J. DASSOW, G. PIGHIZZINI, B. TRUTHE (eds.), *DCFS 2009*. EPTCS 3, 2009, 121–130.
<https://doi.org/10.4204/EPTCS.3.11>
- [5] A. N. MASLOV, Estimates of the number of states of finite automata. *Soviet Math. Doklady* **11** (1970) 5, 1373–1375.
- [6] G. PIGHIZZINI, J. O. SHALLIT, Unary language operations, state complexity and Jacobsthal’s function. *Internat. J. Found. Comput. Sci.* **13** (2002) 1, 145–159.
<https://doi.org/10.1142/S012905410200100X>
- [7] M. SIPSER, *Introduction to the Theory of Computation*. Cengage Learning, 2012.
- [8] S. YU, Q. ZHUANG, K. SALOMAA, The state complexities of some basic operations on regular languages. *Theoret. Comput. Sci.* **125** (1994) 2, 315–328.
[http://dx.doi.org/10.1016/0304-3975\(92\)00011-F](http://dx.doi.org/10.1016/0304-3975(92)00011-F)