



# State Complexity of Simple Splicing

Lila Kari and Timothy Ng<sup>(✉)</sup>

School of Computer Science, University of Waterloo,  
Waterloo, Ontario N2L 3G1, Canada  
{lila.kari,tim.ng}@uwaterloo.ca

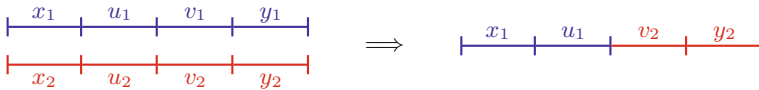
**Abstract.** Splicing, as a binary word/language operation, was inspired by the DNA recombination under the action of restriction enzymes and ligases, and was first introduced by Tom Head in 1987. Splicing systems as generative mechanisms were defined as consisting of an initial starting set of words called an axiom set, and a set of splicing rules—each encoding a splicing operation—, as their computational engine to iteratively generate new strings starting from the axiom set. Since finite splicing systems (splicing systems with a finite axiom set and a finite set of splicing rules) generate a subclass of the family of regular languages, descriptive complexity questions about splicing systems can be answered in terms of the size of the minimal deterministic finite automata that recognize their languages. In this paper we focus on a particular type of splicing systems, called simple splicing systems, where the splicing rules are of a particular form. We prove a tight state complexity bound of  $2^n - 1$  for (semi-)simple splicing systems with a regular initial language with state complexity  $n \geq 3$ . We also show that the state complexity of a (semi-)simple splicing system with a finite initial language is at most  $2^{n-2} + 1$ , and that whether this bound is reachable or not depends on the size of the alphabet and the number of splicing rules.

## 1 Introduction

In [10] Head described a language-theoretic operation, called *splicing*, which models DNA recombination, a cut-and-paste operation on DNA double-stranded molecules, under the action of restriction enzymes and ligases. A *splicing system* is a formal language model which consists of a set of *initial words*,  $I$  (representing double-stranded DNA strings), and a set of *splicing rules*  $R$  (representing restriction enzymes). The most commonly used definition for a splicing rule is a quadruplet of words  $r = (u_1, v_1; u_2, v_2)$ . This rule splices two words  $x_1u_1v_1y_1$  and  $x_2u_2v_2y_2$ : the words are cut between the factors  $u_1, v_1$ , respectively  $u_2, v_2$ , and the prefix (the left segment) of the first word is recombined by catenation with the suffix (the right segment) of the second word; see Fig. 1 and also [16]. The words  $u_1v_1$  and  $u_2v_2$  are the restriction sites in the rule  $r$ . A splicing system generates a language which contains every word that can be obtained by successively applying rules to axioms and the intermediately produced words.

The most natural variant of splicing systems, often referred to as *finite splicing systems*, is to consider a finite set of axioms and a finite set of rules.

Several different types of splicing systems have been proposed in the literature, and Bonizzoni et al. [1] showed that the classes of languages they generate are related. Shortly after the introduction of splicing in formal language theory, Culik II and Harju [4] proved that finite splicing systems can only generate regular languages; see also [11, 15]. Gatterdam [7] gave  $(aa)^*$  as an example of a regular language which cannot be generated by a finite splicing system; thus, the class of languages generated by finite splicing systems is strictly included in the class of regular languages.



**Fig. 1.** Splicing of the words  $x_1u_1v_1y_1$  and  $x_2u_2v_2y_2$  by the rule  $r = (u_1, v_1; u_2, v_2)$ .

Descriptive complexity considers the complexity of a language in terms of the size of a computational device (in this case splicing system) that generates or recognizes it. For instance, Mateescu et al. [14] consider a number of descriptive complexity measures for simple splicing systems, such as the number of rules, the number of words in the initial language, the maximum length of a word in the initial language, and the sum of the lengths of all words in the initial language. Loos et al. [13] consider the descriptive complexity of finite splicing systems by using the number of rules, the length of the rules, and the size of the initial language as complexity measures. Păun [16] proposed using the radius, the largest  $u_i$  in a rule, as a descriptive complexity measure.

As the class of languages generated by splicing systems forms a subclass of the family of regular languages, their descriptive complexity can also be considered in terms of the finite automata that recognize them. For example, Loos et al. [13] gave a bound on the number of states required for a nondeterministic finite automaton to recognize the language generated by an equivalent finite splicing system.

We focus our attention on simple splicing systems, that is, splicing systems where the rules  $(u_1, v_1; u_2, v_2)$  are of a particular form:  $u_1 = u_2 = a$ , are singleton letters, and  $v_1 = v_2 = \varepsilon$  are the empty word. The descriptive complexity of simple splicing systems was considered by Mateescu et al. [14] in terms of the size of a right linear grammar that generates a simple splicing language. Here we consider the descriptive complexity of simple splicing systems in terms of deterministic state complexity [6]. In other words, we measure the descriptive complexity of a simple splicing system in terms of the size of the minimal deterministic finite automaton that recognizes the language generated by the splicing system.

In this paper, we prove tight state complexity bounds for simple and semi-simple splicing systems with regular and finite initial languages. In Sect. 2, we

fix notation and definitions for simple splicing systems. We consider the state complexity of simple splicing systems with regular and finite initial languages in Sect. 3. In Sect. 4, we give tight state complexity bounds for semi-simple splicing systems with finite initial languages. We consider the state complexity of the crossover operation related to simple splicing systems in Sect. 5.

## 2 Preliminaries

Let  $\Sigma$  be a finite alphabet. We denote by  $\Sigma^*$  the set of all finite words over  $\Sigma$ , including the empty word, which we denote by  $\varepsilon$ . We denote the length of a word  $w$  by  $|w| = n$ . If  $w = xyz$  for  $x, y, z \in \Sigma^*$ , we say that  $x$  is a prefix of  $w$ ,  $y$  is a factor of  $w$ , and  $z$  is a suffix of  $w$ .

A deterministic finite automaton (DFA) is a tuple  $A = (Q, \Sigma, \delta, q_0, F)$  where  $Q$  is a finite set of states,  $\Sigma$  is an alphabet,  $\delta$  is a function  $\delta : Q \times \Sigma \rightarrow Q$ ,  $q_0 \in Q$  is the initial state, and  $F \subseteq Q$  is a set of final states. We extend the transition function  $\delta$  to a function  $Q \times \Sigma^* \rightarrow Q$  in the usual way. A DFA  $A$  is complete if  $\delta$  is defined for all  $q \in Q$  and  $a \in \Sigma$ . We will make use of the notation  $q \xrightarrow{w} q'$  for  $\delta(q, w) = q'$ , where  $w \in \Sigma^*$  and  $q, q' \in Q$ . A state  $q \in Q$  is called a sink state if  $\delta(q, a) = q$  for all  $a \in \Sigma$  and  $q \notin F$ .

Each letter  $a \in \Sigma$  defines a transformation of the state set  $Q$ . Let  $\delta_a : Q \rightarrow Q$  be the transformation on  $Q$  induced by  $a$ , defined by  $\delta_a(q) = \delta(q, a)$ . We extend this definition to words by composing the transformations  $\delta_w = \delta_{a_1} \circ \delta_{a_2} \circ \dots \circ \delta_{a_n}$  for  $w = a_1 a_2 \dots a_n$ . We denote by  $\text{im } \delta_a$  the image of  $\delta_a$ , defined  $\text{im } \delta_a = \{\delta(p, a) \mid p \in Q\}$ .

The language recognized or accepted by  $A$  is  $L(A) = \{w \in \Sigma^* \mid \delta(q_0, w) \in F\}$ . A state  $q$  is called *reachable* if there exists a string  $w \in \Sigma^*$  such that  $\delta(q_0, w) = q$ . Two states  $p$  and  $q$  of  $A$  are said to be *equivalent* if  $\delta(p, w) \in F$  if and only if  $\delta(q, w) \in F$  for every word  $w \in \Sigma^*$ . A DFA  $A$  is minimal if each state  $q \in Q$  is reachable from the initial state and no two states are equivalent. The state complexity of a regular language  $L$  is the number of states of the minimal complete DFA recognizing  $L$  [6].

A nondeterministic finite automaton (NFA) is a tuple  $A = (Q, \Sigma, \delta, I, F)$  where  $Q$  is a finite set of states,  $\Sigma$  is an alphabet,  $\delta$  is a function  $\delta : Q \times \Sigma \rightarrow 2^Q$ ,  $I \subseteq Q$  is a set of initial states, and  $F$  is a set of final states. The language recognized by an NFA  $A$  is  $L(A) = \{w \in \Sigma^* \mid \bigcup_{q \in I} \delta(q, w) \cap F \neq \emptyset\}$ . As with DFAs, transitions of  $A$  can be viewed as transformations on the state set. Let  $\delta_a : Q \rightarrow 2^Q$  be the transformation on  $Q$  induced by  $a$ , defined by  $\delta_a(q) = \delta(q, a)$ . The image of  $\delta_a$  is defined by  $\text{im } \delta_a = \{\delta(p, a) \mid p \in Q\}$ . We make use of the notation  $P \xrightarrow{w} P'$  for  $P' = \bigcup_{q \in P} \delta(q, w)$ , where  $w \in \Sigma^*$  and  $P, P' \subseteq Q$ .

### 2.1 Simple Splicing Systems

In this paper we will use the notation of Păun [16], even though simple splicing systems can be defined using any of the three definitions of splicing. The splicing operation is defined via sets of quadruples  $r = (\alpha_1, \alpha_2; \alpha_3, \alpha_4)$

with  $\alpha_1, \alpha_2, \alpha_3, \alpha_4 \in \Sigma^*$  called splicing rules. For two strings  $x = x_1\alpha_1\alpha_2x_2$  and  $y = y_1\alpha_3\alpha_4y_2$ , applying the rule  $r = (\alpha_1, \alpha_2; \alpha_3, \alpha_4)$  produces a string  $z = x_1\alpha_1\alpha_4y_2$ , which we denote by  $(x, y) \vdash^r z$ .

A *splicing scheme* is a pair  $\sigma = (\Sigma, \mathcal{R})$  where  $\Sigma$  is an alphabet and  $\mathcal{R}$  is a set of splicing rules. For a splicing scheme  $\sigma = (\Sigma, \mathcal{R})$  and a language  $L \subseteq \Sigma^*$ , we denote by  $\sigma(L)$  the language

$$\sigma(L) = L \cup \{z \in \Sigma^* \mid (x, y) \vdash^r z, \text{ where } x, y \in L, r \in \mathcal{R}\}.$$

Then we define  $\sigma^0(L) = L$  and  $\sigma^{i+1}(L) = \sigma(\sigma^i(L))$  for  $i \geq 0$  and

$$\sigma^*(L) = \lim_{i \rightarrow \infty} \sigma^i(L) = \bigcup_{i \geq 0} \sigma^i(L).$$

For a splicing scheme  $\sigma = (\Sigma, \mathcal{R})$  and an initial language  $L \subseteq \Sigma^*$ , we say the triple  $H = (\Sigma, \mathcal{R}, L)$  is a *splicing system*. The language generated by  $H$  is defined by  $L(H) = \sigma^*(L)$ .

Mateescu et al. [14] define a restricted class of splicing systems called simple splicing systems. A *simple splicing system* is a triple  $H = (\Sigma, M, I)$ , where  $\Sigma$  is an alphabet,  $M \subseteq \Sigma$  is a set of markers, and  $I$  is a finite initial language over  $\Sigma$ . For  $a \in M$ , we have  $(x, y) \vdash^a z$  if and only if  $x = x_1ax_2$ ,  $y = y_1ay_2$ , and  $z = x_1ay_2$  for some  $x_1, x_2, y_1, y_2 \in \Sigma^*$ .

In other words, a simple splicing system is a system in which the set of rules is  $\mathcal{M} = \{(a, \varepsilon; a, \varepsilon) \mid a \in M\}$  and the initial language  $I$  is finite. Since the rules are determined solely by our choice of  $M \subseteq \Sigma$ , the set  $M$  is used in the definition of the simple splicing system rather than the set of rules  $\mathcal{M}$ . Based on these properties, one can deduce that the class of languages generated by simple splicing systems is subregular [4, 15]. Mateescu et al. [14] show that these languages form a proper subclass of the extended star-free languages.

In this paper, we will relax the condition that the initial language of a simple splicing system must be a finite language. We will consider also simple splicing systems with regular initial languages. By [16], it is clear that such a splicing system will also produce a regular language. In the following, we will use the convention that  $I$  denotes a finite language and  $L$  denotes an infinite language.

### 3 State Complexity of Simple Splicing

In this section, we will give tight state complexity bounds for simple splicing systems with regular and finite initial languages. First, we will define an NFA that recognizes the language of a given simple splicing system. The construction follows a more general construction due to Loos et al. [13] for finite splicing systems. This construction is a simplification of a construction by Pixton [15], which itself is a simplification of the original proof of regularity of finite splicing due to Culik II and Harju [4].

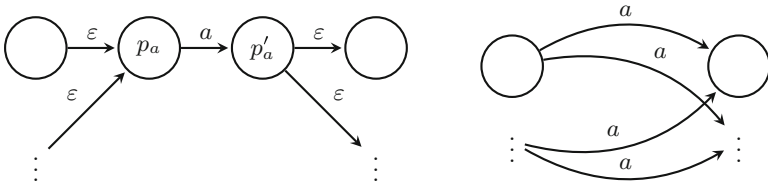
**Proposition 1.** *Let  $H = (\Sigma, M, L)$  be a simple splicing system with a regular initial language  $L$  and let  $L$  be recognized by a DFA with  $n$  states. Then there exists an NFA  $A'_H$  with  $n$  states such that  $L(A'_H) = L(H)$ .*

*Proof.* Let  $H = (\Sigma, M, L)$  and let  $A = (Q, \Sigma, \delta, q_0, F)$  be a DFA for  $L$ . We will define the NFA  $A_H = (Q', \Sigma, \delta', q_0, F)$ , where  $Q' = Q \cup Q_M$  with  $Q_M = \{p_a, p'_a \mid a \in M\}$  and the transition function  $\delta'$  is defined

- $\delta'(q, a) = \{\delta(q, a)\}$  if  $q \in Q$  and  $a \in \Sigma$ ,
- $\delta'(q, \varepsilon) = \{p_a\}$  if  $q \in Q$ ,  $a \in M$ , and  $\delta(q, a)$  is not the sink state,
- $\delta'(p_a, a) = \{p'_a\}$  if  $p_a \in Q_M$ ,
- $\delta'(p'_a, \varepsilon) = \text{im } \delta_a$  if  $p'_a \in Q_M$  and  $a \in M$ .

First, we describe the construction of [13]. Let  $\mathcal{M} = \{(a, \varepsilon; a, \varepsilon) \mid a \in M\}$  be the set of rules for  $H$ . For each rule  $(\alpha_1, \alpha_2; \alpha_3, \alpha_4) \in \mathcal{M}$ , add new states and transitions corresponding to  $\alpha_1\alpha_4$  and  $\alpha_3\alpha_2$ . That is, if  $\alpha_1 = a_1 \cdots a_i$ ,  $\alpha_2 = b_1 \cdots b_j$ ,  $\alpha_3 = c_1 \cdots c_k$ , and  $\alpha_4 = d_1 \cdots d_\ell$ , then add states and transitions corresponding to a path  $r_0 \xrightarrow{a_1} r_1 \xrightarrow{a_2} \cdots \xrightarrow{d_\ell} r_{i+\ell}$  for  $\alpha_1\alpha_4$  and a path  $s_0 \xrightarrow{c_1} s_1 \xrightarrow{c_2} \cdots \xrightarrow{b_j} s_{j+k}$  corresponding to  $\alpha_3\alpha_2$ . Now consider each path  $q \xrightarrow{\alpha_1\alpha_2} q'$  in  $A$  such that  $q$  is reachable from the initial state  $q_0$  and a final state of  $A$  is reachable from  $q'$ . We add an  $\varepsilon$ -transition from  $q$  to  $r_0$  and from  $s_{j+k}$  to  $q'$ . Similarly, for each path  $t \xrightarrow{\alpha_3\alpha_4} t'$ , add  $\varepsilon$ -transitions from  $t$  to  $s_0$  and from  $r_{i+\ell}$  to  $t'$ .

Now, since  $H$  is a simple splicing system, this construction can be simplified further. Since every rule of  $H$  is of the form  $(a, \varepsilon; a, \varepsilon)$ , we only need to add states and transitions for  $p_a \xrightarrow{a} p'_a$  for each rule. Then add  $\varepsilon$ -transitions from states  $q$  of  $A$  to  $p_a$  if  $q$  has an outgoing transition on  $a$  to a non-sink state of  $A$ . From each state  $p'_a$ , add  $\varepsilon$ -transitions to each state of  $A$  with an incoming transition on  $a$ . Recall that  $\text{im } \delta_a$  is the image of the transformation of  $\delta$  induced by  $a$ , and therefore it is the set of states of  $A$  with an incoming transition on  $a$ .



**Fig. 2.** New states and transitions for  $a \in M$  (left), after  $\varepsilon$ -removal (right).

Finally, we can simplify this NFA by removing  $\varepsilon$ -transitions in the usual way to obtain an NFA  $A'_H = (Q, \Sigma, \delta', q_0, F)$ , where

$$\delta'(q, a) = \begin{cases} \{\delta(q, a)\}, & \text{if } \delta(q, a) \text{ is the sink state;} \\ \{\delta(q, a)\}, & \text{if } a \notin M; \\ \text{im } \delta_a, & \text{if } a \in M. \end{cases}$$

Figure 2 illustrates the new states and transitions added for  $a \in M$  before and after  $\varepsilon$ -removal. Observe that by removing the  $\varepsilon$ -transitions, we also remove

the states that were initially added earlier in the construction of  $A_H$ . Thus, the state set of  $A'_H$  is exactly the state set of the DFA  $A$  recognizing  $L$ .  $\square$

Given a splicing system  $H = (\Sigma, M, L)$ , one can obtain a DFA that recognizes  $L(H)$  by performing the subset construction on  $A'_H$ . As shown in Proposition 1, if  $L$  is recognized by a DFA with  $n$  states, then  $A'_H$  also has  $n$  states. By applying the subset construction and observing that the empty set is not reachable from any subset of  $Q$  in  $A'_H$ , this gives an upper bound of  $2^n - 1$  states for a DFA equivalent to  $A'_H$ .

We will now show that there exists a family of regular languages  $L_n$  with state complexity  $n$  such that a simple splicing system  $H = (\Sigma, M, L_n)$  with one marker requires  $2^n - 1$  states for an equivalent DFA to recognize it.

**Proposition 2.** *For  $|\Sigma| \geq 3$  and  $n \geq 3$ , there exists a simple splicing system with a regular initial language  $H = (\Sigma, M, L_n)$  with  $|M| = 1$  where  $L_n$  is a regular language with state complexity  $n$  such that the minimal DFA for  $L(H)$  requires at least  $2^n - 1$  states.*

Proposition 2 is proved via the family of languages  $L_n$  accepted by DFAs  $A_n$ , shown in Fig. 3, with  $M = \{c\}$ .

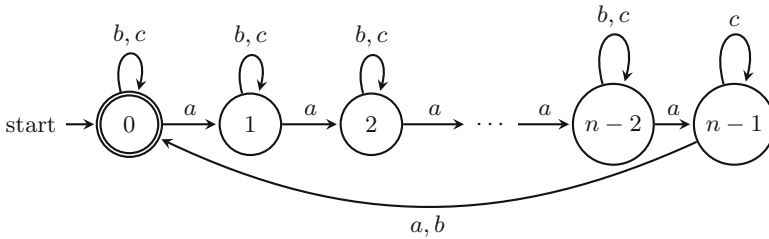


Fig. 3. The DFA  $A_n$ .

Together, Propositions 1 and 2 give the following result.

**Theorem 3.** *For a simple splicing system with a regular initial language  $H = (\Sigma, M, L_n)$  where  $M \subseteq \Sigma$  and  $L_n \subseteq \Sigma^*$  has state complexity  $n$ , the state complexity of  $L(H)$  is at most  $2^n - 1$  and this bound can be reached in the worst case.*

We will now consider simple splicing systems with a finite initial language. We will show that the upper bound of Proposition 1 is not reachable in this case.

**Proposition 4.** *Let  $H = (\Sigma, M, I)$  be a simple splicing system with a finite initial language, where  $I$  is a finite language recognized by a DFA  $A$  with  $n$  states. Then a DFA recognizing  $L(H)$  requires at most  $2^{n-2} + 1$  states.*

We will show that this bound is reachable. We note that the lower bound witness used in the following lemma is defined over an alphabet with size exponential in the number of states of the DFA recognizing the initial language.

**Lemma 5.** *There exists a simple splicing system with a finite initial language  $H = (\Sigma, M, I_n)$  where  $I_n$  is a finite language with state complexity  $n$  such that a DFA recognizing  $L(H)$  requires  $2^{n-2} + 1$  states.*

Together, Proposition 4 and Lemma 5 give the following result.

**Theorem 6.** *For a simple splicing system with a finite initial language  $H = (\Sigma, M, I_n)$  where  $M \subseteq \Sigma$  and  $I_n \subseteq \Sigma^*$  has state complexity  $n$ , the state complexity of  $L(H)$  is at most  $2^{n-2} + 1$  and this bound can be reached in the worst case.*

The bound of Lemma 5 is reached by a witness defined over an alphabet size of  $2^{n-3} + 1$ . An obvious question is whether this bound can be reached via a smaller alphabet. We will consider in the following the state complexity of simple splicing systems with a finite initial language for small, fixed alphabets. We begin with a general observation on the transition function of a DFA recognizing the language of a simple splicing system.

**Lemma 7.** *Let  $H = (\Sigma, M, L)$  be a simple splicing system with a regular initial language and let  $A_H$  be an NFA recognizing  $L(H)$ . If  $a \in M$  and  $\delta'$  is the transition function of  $A_H$ , then  $|\text{im } \delta'_a| = 2$ .*

First, we will consider simple splicing systems with a finite initial language defined over a unary alphabet.

**Proposition 8.** *Let  $H = (\{a\}, M, I)$  be a simple splicing system where  $M$  is nonempty and  $I$  is a finite language containing a word of length at least 2. Then the minimal DFA recognizing  $L(H)$  has exactly two states.*

Next, we consider simple splicing systems with a finite initial language defined over a binary alphabet. We will show that the small size of the alphabet restricts the number of transformations that can be performed on the state set and that the upper bound on the number of states falls far below the upper bound of Proposition 4 as a result.

**Proposition 9.** *Let  $H = (\{a, b\}, M, I)$  be a simple splicing system where  $I$  is a finite language with state complexity  $n$ . Then the minimal DFA recognizing  $L(H)$  has at most  $2n - 3$  states and this bound is reachable in the worst case.*

We will now consider the state complexity of simple splicing systems with a finite initial language defined over a ternary alphabet. We will show that the upper bound of  $2^{n-2} + 1$  from Proposition 4 cannot be reached with an alphabet of size 3.

**Proposition 10.** *Let  $H = (\{a, b, c\}, M, I)$  be a simple splicing system where  $I$  is a finite language with state complexity  $n$ . Then the minimal DFA recognizing  $L(H)$  has at most  $2^{\frac{n}{2}} + 1$  states if  $n$  is even and  $3 \cdot 2^{\frac{n-3}{2}} + 1$  states if  $n$  is odd.*

We note that the upper bound of the previous lemma is similar to the state complexity of the reversal operation on finite languages [2]. We will use this result as inspiration for a family of lower bound witnesses in the following lemma.

**Lemma 11.** *There exists a family of finite languages  $I_n \subseteq \{a, b, c\}^*$ , for  $n \geq 4$ , recognized by a DFA with  $n$  states such that the minimal DFA for a simple splicing system  $H = (\{a, b, c\}, M, I_n)$  requires at least  $2^{\frac{n}{2}} + 1$  states if  $n$  is even and  $3 \cdot 2^{\frac{n-3}{2}} + 1$  states if  $n$  is odd.*

The family of witness languages  $I_n$  used to prove Lemma 11 is accepted by DFAs  $A_n$ , shown in Fig. 4, with  $M = \{c\}$ .

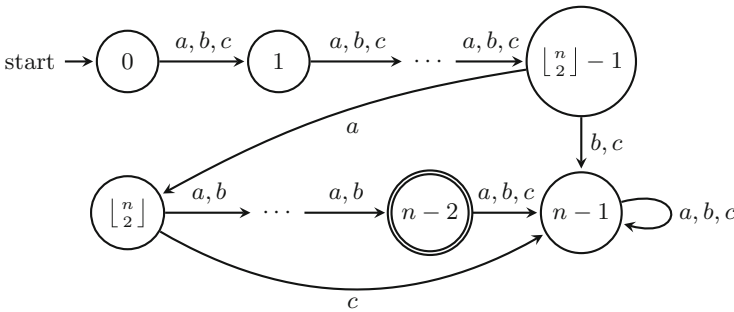


Fig. 4. The ternary witness DFA  $A_n$ .

Together, Proposition 10 and Lemma 11 give us the following theorem.

**Theorem 12.** *For a simple splicing system with a finite initial language  $H = (\Sigma, M, I_n)$  where  $|\Sigma| = 3$ ,  $M \subseteq \Sigma$ , and  $I_n \subseteq \Sigma^*$  has state complexity  $n$ , the state complexity of  $L(H)$  is at most  $2^{\frac{n}{2}} + 1$  states if  $n$  is even and  $3 \cdot 2^{\frac{n-3}{2}} + 1$  states if  $n$  is odd and this bound can be reached in the worst case.*

## 4 State Complexity of Semi-simple Splicing

In this section, we will give tight state complexity bounds for semi-simple splicing systems with regular and finite initial languages. In particular, we will show that the upper bound is reachable for semi-simple splicing systems with a finite initial language defined over a fixed alphabet.

Goode and Pixton [9] generalize simple splicing systems by defining semi-simple splicing systems. A splicing system is semi-simple if every rule is of the form  $(a, \varepsilon; b, \varepsilon)$  for  $a, b \in \Sigma$ . Again, rather than explicitly define a set of rules  $\mathcal{M}$ , we refer instead to the set  $M^{(2)} \subseteq \Sigma \times \Sigma$  of pairs of symbols, which determines the set of rules. As with simple splicing systems, one can conclude that the class of languages generated by semi-simple splicing systems is subregular [4, 15].



In the following, we will give a construction for an NFA that recognizes the language generated by a semi-simple splicing system. As with the NFA for simple splicing systems from Proposition 1, the construction will follow the more general construction for finite splicing systems of Loos et al. [13].

**Proposition 13.** *Let  $H = (\Sigma, M^{(2)}, L)$  be a semi-simple splicing system with a regular initial language. Then there exists an NFA  $B'_H$  with  $n$  states such that  $L(B'_H) = L(H)$ .*

It is clear from Proposition 13 that for a given regular language  $L$ , the language of a semi-simple splicing system  $H = (\Sigma, M^{(2)}, L)$  can require  $2^n - 1$  states in the worst case. Since a simple splicing system is also a semi-simple splicing system, the lower bound witness from Proposition 2 holds. Therefore, we can focus on the more interesting case of semi-simple splicing systems with finite initial languages. First, we observe that even with semi-simple splicing rules, the upper bound on the number of states for a DFA recognizing a semi-simple splicing system with a finite initial language remains the same.

**Proposition 14.** *Let  $H = (\Sigma, M^{(2)}, I)$  be a semi-simple splicing system with a finite initial language where  $I$  is a finite language recognized by a DFA  $A$  with  $n$  states. Then a DFA recognizing  $L(H)$  requires at most  $2^{n-2} + 1$  states.*

The proof of this fact is identical to the proof of Proposition 4.

Recall from Lemma 5, that the lower bound witness for simple splicing systems with a finite initial language was defined over an alphabet with size exponential in the state complexity of the initial language. We will show in the following lemma that for semi-simple splicing systems with a finite initial language, a lower bound witness defined over an alphabet of size 3 exists.

**Lemma 15.** *Let  $n \geq 4$ . Then there exists a semi-simple splicing system with a finite initial language  $H = (\Sigma, M^{(2)}, I_n)$  where  $|\Sigma| = 3$  and  $I_n$  is a finite language with state complexity  $n$  such that  $L(H)$  is recognized by a DFA that requires at least  $2^{n-2} + 1$  states.*

The family of witness languages  $I_n$  of Lemma 15 is accepted by DFAs  $A_n$ , shown in Fig. 5, with  $\Sigma = \{a, b, c\}$  and  $M^{(2)} = \{(a, c)\}$ .

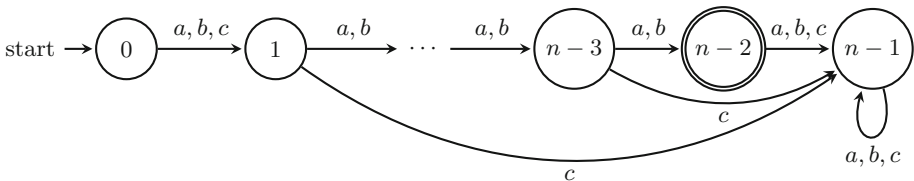


Fig. 5. The ternary witness DFA  $A_n$ .

From Proposition 14 and Lemma 15, we have the following result.

**Theorem 16.** *For a semi-simple splicing system with a finite initial language  $H = (\Sigma, M^{(2)}, I_n)$  where  $M \subseteq \Sigma$  and  $I_n \subseteq \Sigma^*$  has state complexity  $n$ , the state complexity of  $L(H)$  is at most  $2^{n-2} + 1$  and this bound can be reached in the worst case.*

## 5 State Complexity of the Crossover Operation

In this section, we will give tight state complexity bounds for the crossover operation [3], which can be thought of as a single step of semi-simple splicing. Mateescu et al. [14] gave an algebraic characterization of the class of languages generated by simple splicing systems based on the crossover operation therein. A similar such characterization for the class of languages generated by semi-simple splicing systems is given by Ceterchi [3].

For  $M = M_1 \times M_2 \subseteq \Sigma \times \Sigma$ , define the operation  $\diamond_M$  on two strings  $u, v \in \Sigma^+$  by

$$u \diamond_M v = \begin{cases} u'av', & \text{if } u = u'a, v = bv' \text{ for } (a, b) \in M, u', v' \in \Sigma^*; \\ \text{undefined,} & \text{otherwise.} \end{cases}$$

Then for two languages  $L_1, L_2 \subseteq \Sigma^*$ , we have

$$L_1 \diamond_M L_2 = \{x \diamond_M y \mid x \in L_1, y \in L_2\}.$$

The operation  $\diamond_M$  is a variant of the Latin product defined in [8]. Based on  $\diamond_M$ , we define the *crossover operation*  $\sharp_M$  for  $M \subseteq \Sigma \times \Sigma$  and two languages  $L_1, L_2 \subseteq \Sigma^*$  by

$$L_1 \sharp_M L_2 = \text{pref}(L_1) \diamond_M \text{suff}(L_2),$$

where  $\text{pref}(L_1)$  is the set of prefixes of words in  $L_1$  and  $\text{suff}(L_2)$  is the set of suffixes of words in  $L_2$ . From this definition, the operation  $\sharp_M$  can be viewed as a combination of operations under each of which the regular languages are closed. Therefore, it is easy to see that the regular languages are closed under  $\sharp_M$ .

Note that by restricting  $M$  to pairs  $(a, a)$  for  $a \in \Sigma$ , we get an operation that can be thought of as a single step of simple splicing. The operation  $\sharp_M$ , when restricted to pairs of the form  $(a, a)$  has some similarities to many operations that have been studied in the literature, such as the chop operation [12] and the word blending operation [5]. In fact, word blending can be seen as a special case of the crossover operation, taking  $M = \{(a, a) \mid a \in \Sigma\}$ .

We will now give a DFA construction for the crossover of two regular languages.

**Proposition 17.** *Let  $A$  and  $B$  be two DFAs defined over  $\Sigma$  with  $m$  and  $n$  states, respectively. Then for any  $M \subseteq \Sigma \times \Sigma$ , there exists a DFA  $C$  such that  $L(C) = L(A) \sharp_M L(B)$  and  $C$  has at most  $m \cdot 2^n$  states.*

*Proof.* Let  $A = (Q_A, \Sigma, \delta_A, s_A, F_A)$  and  $B = (Q_B, \Sigma, \delta_B, s_B, F_B)$  be two DFAs. We will construct a DFA  $C = (Q_C, \Sigma, \delta_C, s_C, F_C)$  that recognizes  $A\#_M B$  for some  $M \subseteq \Sigma \times \Sigma$ , defined by

- $Q_C = Q_A \times 2^{Q_B}$ ,
- $s_C = \langle s_A, \emptyset \rangle$ ,
- $F_C = \{ \langle q, P \rangle \in Q_A \times 2^{Q_B} \mid P \cap F_B \neq \emptyset \}$ ,

and the transition function  $\delta_C$  is defined for  $q \in Q_A$ ,  $P \subseteq Q_B$ , and  $a \in \Sigma$  by  $\delta_C(\langle q, P \rangle, a) = \langle q', P' \rangle$ , where  $q' = \delta_A(q, a)$  and

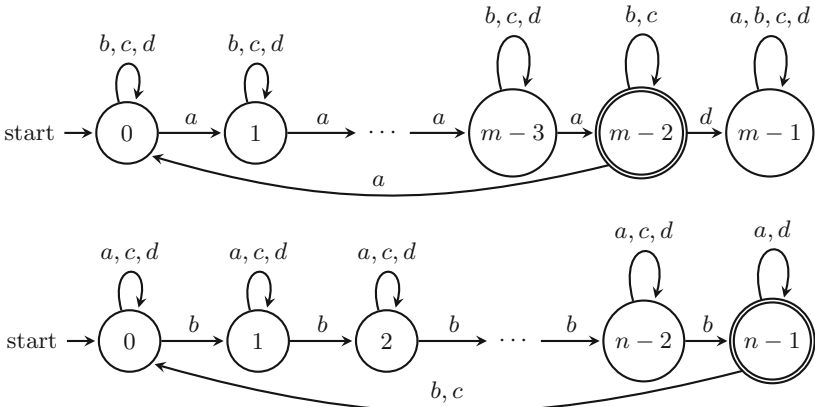
$$P' = \begin{cases} \text{im}(\delta_B)_b, & \text{if } (a, b) \in M \text{ and } q' \text{ is not a sink state;} \\ \bigcup_{p \in P} \delta_B(p, a), & \text{otherwise.} \end{cases}$$

Informally, the machine traces a computation of  $A$  and computations of  $B$ . For each pair  $(a, b) \in M$ , whenever a transition on  $a$  occurs in  $A$ , all states of  $B$  with incoming transitions on  $b$  are added to the computation. It is clear from the definition of  $C$  that  $L(C) = L(A)\#_M L(B)$  and it has at most  $m \cdot 2^n$  states. □

We will show that the bound of Proposition 17 is reachable, even when  $M$  is restricted to pairs of the form  $(a, a)$ .

**Lemma 18.** *There exist languages  $A_m$  and  $B_n$  over  $\Sigma$  with  $|\Sigma| \geq 4$  recognized by DFAs with  $m$  and  $n$  states, respectively, and a subset  $M \subseteq \Sigma \times \Sigma$  such that the minimal DFA for  $L(A_m)\#_M L(B_n)$  requires at least  $m \cdot 2^n$  states.*

The families of witness languages of Lemma 18 are accepted by DFAs  $A_m$  and  $B_n$ , shown in Fig. 6, with  $M = \{(d, d)\}$ .



**Fig. 6.** The DFAs  $A_m$  (above) and  $B_n$  (below).

Together, Proposition 17 and Lemma 18 give us the following theorem.

**Theorem 19.** *For regular languages  $L_m$  and  $L_n$ , with  $m, n \geq 3$ , defined over an alphabet  $\Sigma$ , with  $|\Sigma| \geq 4$ , and a subset  $M \subseteq \Sigma \times \Sigma$ , if  $L_m$  has state complexity  $m$  and  $L_n$  has state complexity  $n$ , then  $L_m \sharp_M L_n$  has state complexity at most  $m \cdot 2^n$  and this bound can be reached in the worst case.*

## 6 Conclusion

We have given tight bounds for the state complexity of simple and semi-simple splicing systems and the associated crossover operation. In almost all cases, the exponential upper bound was easily reached via splicing systems defined over a fixed-size alphabet with one rule. The exception is with simple splicing systems with a finite initial language, where a natural open problem to consider is the worst-case state complexity when the initial languages are defined over alphabets of size between 3 and  $2^{n-3}$ .

## References

1. Bonizzoni, P., Ferretti, C., Mauri, G., Zizza, R.: Separating some splicing models. *Inform. Process. Lett.* **79**(6), 255–259 (2001)
2. Câmpeanu, C., Culik, K., Salomaa, K., Yu, S.: State Complexity of Basic Operations on Finite Languages. In: Boldt, O., Jürgensen, H. (eds.) WIA 1999. LNCS, vol. 2214, pp. 60–70. Springer, Heidelberg (2001). [https://doi.org/10.1007/3-540-45526-4\\_6](https://doi.org/10.1007/3-540-45526-4_6)
3. Ceterchi, R.: An algebraic characterization of semi-simple splicing. *Fund. Inform.* **73**(1–2), 19–25 (2006)
4. Culik II, K., Harju, T.: Splicing semigroups of dominoes and DNA. *Discrete Appl. Math.* **31**(3), 261–277 (1991). [https://doi.org/10.1016/0166-218X\(91\)90054-Z](https://doi.org/10.1016/0166-218X(91)90054-Z)
5. Enaganti, S.K., Kari, L., Ng, T., Wang, Z.: Word blending in formal languages: the brangelina effect. In: Stepney, S., Verlan, S. (eds.) UCNC 2018. LNCS, vol. 10867, pp. 72–85. Springer, Cham (2018). [https://doi.org/10.1007/978-3-319-92435-9\\_6](https://doi.org/10.1007/978-3-319-92435-9_6)
6. Gao, Y., Moreira, N., Reis, R., Yu, S.: A survey on operational state complexity. *J. Autom. Lang. Comb.* **21**(4), 251–310 (2017). <https://doi.org/10.25596/jalc-2016-251>
7. Gatterdam, R.W.: Splicing systems and regularity. *Int. J. Comput. Math.* **31**(1–2), 63–67 (1989). <https://doi.org/10.1080/00207168908803788>
8. Golan, J.S.: *Semirings and their Applications*. Springer, Dordrecht (1999). <https://doi.org/10.1007/978-94-015-9333-5>
9. Goode, E., Pixton, D.: Semi-simple splicing systems. *Where Mathematics. Computer Science, Linguistics and Biology Meet*, pp. 343–352. Springer, Dordrecht (2001). [https://doi.org/10.1007/978-94-015-9634-3\\_30](https://doi.org/10.1007/978-94-015-9634-3_30)
10. Head, T.: Formal language theory and DNA: an analysis of the generative capacity of specific recombinant behaviors. *Bull. Math. Biol.* **49**(6), 737–759 (1987)
11. Head, T., Pixton, D.: Splicing and regularity. In: Esik, Z., Martín-Vide, C., Mitrana, V. (eds.) *Recent Advances in Formal Languages and Applications*. SCI, vol. 25, pp. 119–147. Springer, Heidelberg (2006). [https://doi.org/10.1007/978-3-540-33461-3\\_5](https://doi.org/10.1007/978-3-540-33461-3_5)

12. Holzer, M., Jakobi, S.: Chop operations and expressions: descriptonal complexity considerations. In: Mauri, G., Leporati, A. (eds.) DLT 2011. LNCS, vol. 6795, pp. 264–275. Springer, Heidelberg (2011). [https://doi.org/10.1007/978-3-642-22321-1\\_23](https://doi.org/10.1007/978-3-642-22321-1_23)
13. Loos, R., Malcher, A., Wotschke, D.: Descriptonal complexity of splicing systems. *Internat. J. Found. Comput. Sci.* **19**(4), 813–826 (2008)
14. Mateescu, A., Păun, G., Rozenberg, G., Salomaa, A.: Simple splicing systems. *Discrete Appl. Math.* **84**(1–3), 145–163 (1998). [https://doi.org/10.1016/S0166-218X\(98\)00002-X](https://doi.org/10.1016/S0166-218X(98)00002-X)
15. Pixton, D.: Regularity of splicing languages. *Discrete Appl. Math.* **69**(1–2), 101–124 (1996). [https://doi.org/10.1016/0166-218X\(95\)00079-7](https://doi.org/10.1016/0166-218X(95)00079-7)
16. Păun, G.: On the splicing operation. *Discrete Appl. Math.* **70**(1), 57–79 (1996). [https://doi.org/10.1016/0166-218X\(96\)00101-1](https://doi.org/10.1016/0166-218X(96)00101-1)