

Word blending in formal languages: The Brangelina effect

Srujan Kumar Enaganti¹, Lila Kari², Timothy Ng², and Zihao Wang¹

¹ Department of Computer Science, The University of Western Ontario
London, Ontario N6A 3K7, Canada

`srujankumar@gmail.com`, `zwang688@uwo.ca`

² School of Computer Science, University of Waterloo

Waterloo, Ontario N2L 3GL, Canada

`{lila,tim.ng}@uwaterloo.ca`

Abstract. In this paper we define and investigate a binary word operation that formalizes an experimentally observed outcome of DNA computations, performed to generate a small gene library and implemented using a DNA recombination technique called Cross-pairing Polymerase Chain Reaction (XPCR). The *word blending* between two words xy_1 and y_2wz that share a non-empty overlap w , results in xwz . We study closure properties of families in the Chomsky hierarchy under word blending, language equations involving this operation, and its descriptive state complexity when applied to regular languages. Interestingly, this phenomenon has been observed independently in linguistics, under the name “blend word” or “portmanteau”, and is responsible for the creation of words in the English language such as smog (*smoke* + *fog*), labradoodle (*labrador* + *poodle*), and Brangelina (*Brad* + *Angelina*).

1 Introduction

Cross-pairing Polymerase Chain Reaction (XPCR) is an experimental DNA protocol introduced in [11] for extracting, from a heterogeneous pool of DNA strands, all the strands containing a given substrand. XPCR was then employed to implement several DNA recombination algorithms [13], for the creation of the solution space for a SAT problem [9], and for mutagenesis [12]. The combinatorial power of such a technique has been explained by logical-symbolic schemes in [23], while algorithms to create combinatorial libraries were improved and experimented in [12], [10].

The formal language operation called *overlap assembly*, introduced in [5] under the name of self-assembly, and further investigated in [7, 8, 3], also models a special case of XPCR: The overlap assembly of two strings αx and $x\beta$ that share a non-empty overlap x , results in the string $\alpha x\beta$. A particular case of overlap assembly, called “chop operation”, where the overlap consists of a single letter, was studied in [18, 19], and generalized to an arbitrary length overlap in [20]. Other similar operations have been studied in the literature, such as the “short concatenation” [4], which uses only the maximum-length (possibly

empty) overlap y between operands, the “Latin product” of words [14] where the overlap consists of only one letter, and the operation \otimes which imposes the restriction that the non-overlapping part x is not empty [21]. Overlap assembly can also be considered as a particular case of “semantic shuffle on trajectories” with trajectory $0^*\sigma^+1^*$ or as a generalization of the operation \odot_N from [6] which imposes the length of the overlap to be at least N . Many similar biological phenomena and operations can also be modelled using splicing systems [26, 27]. However, modeling these operations often does not require the full power of splicing. Properties of splicing languages under restrictions such as symmetry and reflexivity have been studied in [2, 15].

Returning to the biological process that motivated the study of overlap assembly, the XPCR procedure has been successfully used to join two different genes if they are attached to compatible primers [10]. Formally, $\alpha A\gamma$ and $\gamma D\beta$ were combined to produce $\alpha A\gamma D\beta$ (here A and D are gene sequences and α , γ and β are primers used). However, when $A = D$, that is, when two sequences containing the same gene were combined by XPCR, the result was not as expected. More specifically, when using XPCR with two strings $\alpha A\gamma$ and $\gamma A\beta$, instead of obtaining the expected $\alpha A\gamma A\beta$, the experiments repeatedly produced the result $\alpha A\beta$.

In this paper, we define and investigate a formal language operation called *word blending*, that formalizes this experimentally observed outcome of XPCR: The *word blending* of two words xAy_1 and y_2Az that share a non-empty overlap A results in xAz . Interestingly, this phenomenon has been observed independently in linguistics [16], under the name “blend word” or “portmanteau”, and is responsible for the creation of words in the English language such as *smog* (*smoke* + *fog*), *labradoodle* (*labrador* + *poodle*), *emoticon* (*emotion* + *icon*), and *Brangelina* (*Brad* + *Angelina*).

The paper is organized as follows. Section 2 details the biological motivation behind the study of word blending, and introduces the main definitions and notations. Section 3 studies closure properties of the families in the Chomsky hierarchy under word blending, its right and left inverses, as well as iterated word blending. Section 4 investigates the decidability of existence of solutions to some language equations involving word blending, and Section 5 studies the descriptive state complexity of this operation when applied to regular languages.

2 Preliminaries

An alphabet Σ is a finite non-empty set of symbols. Σ^* denotes the set of all words over Σ , including the empty word λ , and Σ^+ denotes the set of all non-empty words over Σ . The length of the word w is denoted $lg(w)$. For words $w, x, y, z \in \Sigma^*$ such that $w = xyz$ we call the subwords x , y , and z *prefix*, *infix*, and *suffix* of w , respectively. The sets $\text{pref}(w)$, $\text{inf}(w)$, and $\text{suff}(w)$ contain, respectively, all prefixes, infixes, and suffixes of w . This notation is extended to languages as $\text{suff}(L) = \bigcup_{w \in L} \text{suff}(w)$. The mirror image of a word $w \in \Sigma^*$ is defined as $\text{mi}(\lambda) = \lambda$, and $\text{mi}(w) = a_k \dots a_2 a_1$ if $w = a_1 a_2 \dots a_k$. The definition is extended

to languages in the natural way, by $\text{mi}(L) = \bigcup_{w \in L} \text{mi}(w)$. The complement of a language $L \subseteq \Sigma^*$ is $L^c = \Sigma^* \setminus L$. For two languages L_1 and L_2 , the right quotient of L_1 by L_2 is defined as $L_1 L_2^{-1} = \{u \in \Sigma^* \mid \exists uv \in L_1, v \in L_2\}$, and the left quotient of L_1 by L_2 is defined as $L_2^{-1} L_1 = \{v \in \Sigma^* \mid \exists uv \in L_1, u \in L_2\}$.

The biological phenomenon we model in this paper was observed during the XPCR-based experiments, initially intended to achieve the catenation of two or more genes (genomic DNA strands). It was namely observed in [10] that, in the particular case where the two genes to be catenated were one and the same, that is, when the two input DNA strands were $\alpha A \gamma$ and $\gamma A \beta$ (here A represents a gene sequence), the output of a PCR-based amplification with primers α and β was $\alpha A \beta$. This output was different from the expected $\alpha A \gamma A \beta$, which had been the anticipated result. (Indeed, experiments using XPCR for the purpose of catenating two different genes A and D flanked by primers, that is, when the two input strands were $\alpha A \gamma$ and $\gamma D \beta$, had resulted in the output $\alpha A \gamma D \beta$. This “expected” output of XPCR was modelled by the previously mentioned operation of overlap assembly, given by $\alpha A \gamma + \gamma D \beta = \alpha A \gamma D \beta$.)

Generalizing this experimentally newly-observed phenomenon to the case where the end words of the input strings are different, we model this string recombination as follows. Given two non-empty words x, y over an alphabet Σ , we define the *word blending*, or simply *blending*, of x with y as

$$x \bowtie y = \{z \in \Sigma^+ \mid \exists \alpha, \beta, \gamma_1, \gamma_2 \in \Sigma^*, \exists w \in \Sigma^+ : x = \alpha w \gamma_1, y = \gamma_2 w \beta, z = \alpha w \beta\}.$$

The definition of blending can be extended to languages L_1 and L_2 by

$$L_1 \bowtie L_2 = \bigcup_{x \in L_1, y \in L_2} x \bowtie y.$$

Note that, for a realistic model, we would need additional restrictions such as the fact that the w , γ_1 and γ_2 should be of a sufficient length and should not appear as a substring in the other strings involved.

We can also extend the blending operation to an iterated version on a language. Let $L \subseteq \Sigma^*$ be a language. We define the *iterated (word) blending* of L by $L^{\bowtie_0} = L$ and $L^{\bowtie_i} = L \bowtie L^{\bowtie_{i-1}}$. We define the iterated blending closure of L by

$$L^{\bowtie_*} = \bigcup_{i \geq 0} L^{\bowtie_i}.$$

We observe that the result of the iterated blending operation can be generated by a splicing system with null context splicing rules [17]. Splicing rules in [17] are of the form $(u_1, z, u_2; u_3, z, u_4)$. For such a rule, if we have strings $x = x_1 u_1 z u_2 x_2$ and $y = y_1 u_3 z u_4 y_2$, we obtain the word $x_1 u_1 z u_4 y_2$. A splicing rule is a null context rule when $u_1, u_2, u_3, u_4 = \lambda$. It is easy to see that the language L^{\bowtie_*} can be generated from a splicing scheme with rules of the form $(\lambda, w, \lambda; \lambda, w, \lambda)$ for every word $w \in \Sigma^+$. The relationship between iterated blending and splicing will be discussed in greater detail in Section 3.

3 Closure Properties

In this section, we prove that the families of regular, context-free and recursively enumerable languages are closed under blending, and that the family of context-sensitive languages is not. The section also contains closure properties of Chomsky hierarchy families under the right and left inverse of word blending, as well as under iterated word blending.

The following lemma shows that word-blending is equivalent to a restricted version where only one-letter overlaps are utilized.

Lemma 1. *If x, y are non-empty words in Σ^+ , then*

$$x \bowtie y = \{z \in \Sigma^+ \mid \exists \alpha, \beta, \gamma_1, \gamma_2 \in \Sigma^*, \exists a \in \Sigma : x = \alpha a \gamma_1, y = \gamma_2 a \beta, z = \alpha a \beta\}.$$

This result can be extended to languages in the natural way. Then from this lemma, we can show that the word blending of two languages can be obtained by combining the right quotient, catenation, left quotient and union operations, as follows.

Proposition 2. *Given languages $L_1, L_2 \subseteq \Sigma^+$,*

$$L_1 \bowtie L_2 = \bigcup_{a \in \Sigma} (L_1 (a \Sigma^*)^{-1}) a ((\Sigma^* a)^{-1} L_2).$$

Corollary 3. *Every full AFL is closed under word blending.*

We note that the families of regular languages, context-free languages and recursively enumerable languages are all full AFLs [28].

Proposition 4. *The family of context-sensitive languages is not closed under word blending.*

Proof. Let L_0 be a recursively enumerable language over Σ , that is not context-sensitive. It is known that a context-sensitive language L_1 over $\Sigma \cup \{a, b\}$ with $a, b \notin \Sigma$, can be constructed such that L_1 consists of words of the form Pba^i where $i \geq 0$ and $P \in L_0$ and, in addition, for every $P \in L_0$ there is an $i \geq 0$ such that $Pba^i \in L_1$ (see, e.g., [28]).

Since it is obvious that $L_1 \bowtie \{b\} = \{Pb \mid P \in L_0\}$, which is not context sensitive, it follows that the family of context sensitive languages is not closed under word blending with singleton words. \square

Recall that, given a binary word operation \diamond , the binary word operation \square is called the *right-inverse of \diamond* [22] if and only if for every triplet of words $u, y, w \in \Sigma^*$ the following relation holds: $w \in (u \diamond y)$ if and only if $y \in (u \square w)$. In other words, the operation \square is called the right-inverse of \diamond if it can be used to recover the right operand y in $u \diamond y$, from the other operand u and a word $w \in (u \diamond y)$ in the result. Define now the binary word operation \bowtie^r as $u \bowtie^r w = \bigcup_{a \in \Sigma} \Sigma^* a \left((u (a \Sigma^*)^{-1} a)^{-1} w \right)$. Informally, given a word $w = \alpha a \beta \in (\alpha a \gamma_1 \bowtie \gamma_2 a \beta)$, the operation \bowtie^r outputs the right operand $y = \gamma_2 a \beta$ of word blending, if it is given as inputs the result $w = \alpha a \beta \in (u \bowtie y)$ and the left operand $u = \alpha a \gamma_1$. The definition of \bowtie^r can be extended to languages naturally.

Proposition 5. *The operation \bowtie^r is the right-inverse of \bowtie .*

Proof. If $w \in u \bowtie y$, there exist $\alpha, \beta, \gamma_1, \gamma_2 \in \Sigma^*, b \in \Sigma$ such that $w = \alpha b \beta, u = \alpha b \gamma_1, y = \gamma_2 b \beta$ by Lemma 1. Then, we have that $y = \gamma_2 b \beta \in \Sigma^* b \beta = \Sigma^* b ((\alpha b)^{-1} (\alpha b \beta)) \subseteq \Sigma^* b \left((((\alpha b \gamma_1)(b \Sigma^*)^{-1}) b)^{-1} (\alpha b \beta) \right) \subseteq \bigcup_{a \in \Sigma} \Sigma^* a \left((((\alpha b \gamma_1)(a \Sigma^*)^{-1}) a)^{-1} (\alpha b \beta) \right)$.

If $y \in u \bowtie^r w = \bigcup_{a \in \Sigma} \Sigma^* a \left(((u(a \Sigma^*)^{-1}) a)^{-1} w \right)$, then there exist $b \in \Sigma$, and $\gamma_2 \in \Sigma^*, \gamma_3 \in (u(b \Sigma^*)^{-1}) b$ such that $y = \gamma_2 b (\gamma_3^{-1} w)$. This implies that $w \in (u(b \Sigma^*)^{-1}) b (\gamma_3^{-1} w) = (u(b \Sigma^*)^{-1}) b ((\gamma_2 b)^{-1} (\gamma_2 b (\gamma_3^{-1} w)))$ which is included in $(u(b \Sigma^*)^{-1}) b ((\Sigma^* b)^{-1} y) \subseteq \bigcup_{a \in \Sigma} (u(a \Sigma^*)^{-1}) a ((\Sigma^* a)^{-1} y) = u \bowtie y$. \square

Corollary 6. *The families of regular languages and recursively enumerable languages are closed under the right inverse of the blending. Moreover, if L_1 is an arbitrary language and L_2 is a regular language, then $L_1 \bowtie^r L_2$ is regular; if L_1 is a regular language and L_2 is a context-free language, then $L_1 \bowtie^r L_2$ is context-free.*

Proposition 7. *The family of context-free languages is not closed under the right inverse of blending.*

Proof. Consider the context-free languages $L_1 = \{a \$ (b^{i_1} a^{i_1} \$) \dots (b^{i_n} a^{i_n} \$) \mid n \geq 1, i_m \geq 1 \text{ for } 1 \leq m \leq n\}$, $L_2 = \{(a^{j_1} \$ b^{2j_1}) \dots (a^{j_k} \$ b^{2j_k}) (a^j \$ c^{2j}) \mid j \geq 1, k \geq 1, j_m \geq 1 \text{ for } 1 \leq m \leq k\}$ and the regular language $R = \{\$ c^{2^n}\}$.

We now show that $(L_1 \bowtie^r L_2) \cap R = \{\$ c^{2^n} \mid n \geq 2\}$. Since words in R start with $\$$ and contain only one symbol $\$$, the only cases in which the words in $L_1 \bowtie^r L_2$ have the pattern of the words in R are the cases of word pairs where the overlap letter is $\$$, and a prefix ending in $\$$ in the word from L_1 matches the prefix ending in the last occurrence of $\$$ in the word from L_2 . More precisely, let $u = a \$ b^{i_1} a^{i_1} \$ b^{i_2} a^{i_2} \$ \dots b^{i_m} a^{i_m} \$ \dots b^{i_n} a^{i_n} \$ \in L_1$ and $v = a^{j_1} \$ b^{2j_1} a^{j_2} \$ b^{2j_2} \dots a^{j_m} \$ b^{2j_m} a^j \$ c^{2j} \in L_2$. For a word $w \in (L_1 \bowtie^r L_2)$ to belong to R , we must have

$$a \$ b^{i_1} a^{i_1} \$ b^{i_2} a^{i_2} \$ \dots b^{i_m} a^{i_m} \$ = a^{j_1} \$ b^{2j_1} a^{j_2} \$ b^{2j_2} \dots a^{j_m} \$ b^{2j_m} a^j \$,$$

which implies $j_1 = 1, j_2 = i_1 = 2j_1 = 2, \dots, j = i_m = 2j_m = 2^m$. Thus, $w = \$ c^{2^j} = \$ c^{2^{m+1}}$, which implies $(L_1 \bowtie^r L_2) \cap R = \{\$ c^{2^n} \mid n \geq 2\}$.

Since the family of context-free languages is closed under intersection with regular languages, it follows that it is not closed under the right inverse of blending. \square

Proposition 8. *The family of context-sensitive languages is not closed under the right inverse of blending.*

Recall that given a binary word operation \diamond , the binary word operation \square is called the *left-inverse of \diamond* iff for every triplet of words $x, v, w \in \Sigma^*$ the following relation holds: $w \in (x \diamond v)$ if and only if $x \in (w \square v)$ [22].

Proposition 9. *The left inverse of blending can be expressed using the right inverse of blending, and mirror image as $w \bowtie^l v = \text{mi}(\text{mi}(v) \bowtie^r \text{mi}(w))$.*

Because all families of languages in the Chomsky hierarchy are closed under mirror image, their closure properties under the left-inverse of word blending are the same as their closure properties under the right-inverse of word blending.

We now consider the iterated blending operation \bowtie_* . Recall that, as mentioned in Section 2, for any language $L \subseteq \Sigma^*$, the language $L^{\bowtie*}$ can be generated by a splicing system with null-context splicing rules defined as 6-tuples, as in [17]. As shown in [1], every splicing system where the rules are defined by 6-tuples, can also be implemented by a splicing system as defined in [27], which uses 4-tuple rules (see Definition 10). This connection, together with Proposition 2, allows us to express iterated word blending using so-called simple splicing systems [24], themselves a particular case of splicing systems based on 4-tuple splicing rules.

Definition 10 ([27]). *Let $\sigma = (\Sigma, R)$ be a splicing scheme, where Σ is the alphabet and R is a set of rules $R \subseteq \Sigma^* \# \Sigma^* \$ \Sigma^* \# \Sigma^*$. A rule $(u_1, u_2; u_3, u_4)$ is a word $u_1 \# u_2 \$ u_3 \# u_4 \in R$. For two strings $x, y \in \Sigma^*$, we have*

$$\sigma(x, y) = \{x_1 u_1 u_4 y_2 \mid x = x_1 u_1 u_2 x_2, y = y_1 u_3 u_4 y_2; \\ x_1, x_2, y_1, y_2 \in \Sigma^*, u_1 \# u_2 \$ u_3 \# u_4 \in R\}.$$

For a language L , we define $\sigma(L) = L \cup \bigcup_{x, y \in L} \sigma(x, y)$ and we define the iterated splicing of L by $\sigma^*(L) = \bigcup_{i \geq 0} \sigma^i(L)$ with $\sigma^0(L) = L$ and $\sigma^{i+1}(L) = \sigma(\sigma^i(L))$.

Simple splicing schemes are splicing schemes as above, but restricted to rules of the form $(a, \lambda; a, \lambda)$ for $a \in \Sigma$. Note that for two languages L_1 and L_2 over Σ , we now have that

$$L_1 \bowtie L_2 = \bigcup_{x \in L_1, y \in L_2} \sigma_{\bowtie}(x, y),$$

where σ_{\bowtie} is the simple splicing scheme $\sigma_{\bowtie} = (\Sigma, R)$ with $R = \Sigma \# \lambda \$ \Sigma \# \lambda$. This observation together with Proposition 2 which showed that the word blending of two languages can be written $L_1 \bowtie L_2 = \bigcup_{a \in \Sigma} (L_1 (a \Sigma^*)^{-1}) a ((\Sigma^* a)^{-1} L_2)$, gives us the following result.

Proposition 11. *For any language $L \subseteq \Sigma^*$, we have $\sigma_{\bowtie}(L) = L \bowtie L$ and $\sigma_{\bowtie}^*(L) = L^{\bowtie*}$.*

We note that the splicing scheme σ_{\bowtie} is finite, since the number of rules depends only on the number of symbols in Σ , and it is unary, since the rules use words of length at most 1. We also note that, even though in [24] consideration is restricted to the case when L is a finite language, the properties of the splicing systems obtained therein imply the following closure properties.

Proposition 12. *Every full AFL is closed under iterated word blending.*

Proof. Recall that $L^{\bowtie*} = \sigma_{\bowtie}^*(L)$ and that σ_{\bowtie}^* is finite and unary. For a splicing rule $u_1 \# u_2 \$ u_3 \# u_4$, the words u_1 and u_4 are called visible sites and u_2 and u_3

are invisible sites. In [26], it is shown that full AFLs are closed under regular splicing systems with finitely many visible sites. Since σ_{\bowtie}^* is finite, the rules of σ_{\bowtie}^* contain only finitely many visible sites. \square

Now, we will give an explicit construction for $L^{\bowtie*}$ when L is a regular language. We will require the following lemma concerning the structure of words generated by the iterated blending operation.

Lemma 13. *Let $L \subseteq \Sigma^+$ be a language. Then for each word $w \in L^{\bowtie*}$, there exists $n \in \mathbb{N}$ such that there are words $u_i \in \text{inf}(L)$, $1 \leq i \leq n$ and $\alpha_j \in \Sigma^*$, $1 \leq j \leq n$ and symbols $a_k \in \Sigma$, $1 \leq k \leq n-1$ where*

1. for $n > 1$,
 - (a) $w = \alpha_1 a_1 \alpha_2 a_2 \cdots a_{n-1} \alpha_n$,
 - (b) $u_i = a_{i-1} \alpha_i a_i \in \text{inf}(L)$ for all $2 \leq i \leq n-1$,
 - (c) $u_1 = \alpha_1 a_1 \in \text{pref}(L)$ and $u_n = a_{n-1} \alpha_n \in \text{suff}(L)$,
2. $u_1 = w \in L$ for $n = 1$.

Proposition 14. *Given an NFA A , there exists an NFA A' recognizing the language $L(A)^{\bowtie*}$ which is effectively constructible.*

This construction gives us a way to test whether a regular language L is closed under iterated blending.

Proposition 15. *Let L be a regular language. It is decidable whether or not L is closed under \bowtie_* .*

Let $L, B \subseteq \Sigma^*$ be two languages. We say that B is a base of L (with respect to \bowtie) if $L = B^{\bowtie*}$. In [24], it is shown that it is decidable whether or not a regular language is generated by a simple splicing scheme and a finite language base. Here, we extend the result to consider the case when the base need not be finite.

Theorem 16. *It is decidable whether or not a regular language has a base over \bowtie_* .*

As a consequence, we are able to not only decide whether a regular language is closed under \bowtie_* , but if it is, we know there always exists a finite base that generates it.

Corollary 17. *Let L be a regular language closed under \bowtie_* . Then L can be generated by a finite base.*

Note that in [24] languages generated by simple splicing schemes are assumed to have finite bases by definition. There it was also shown that the class of languages generated by these simple splicing schemes is a subclass of the family of regular languages. Here we do not have the finite base restriction, and Corollary 17 shows that allowing regular bases does not give simple splicing schemes and iterated word blending any more power than restricting bases to be finite.

4 Decision Problems

This section investigates the existence of solutions to language equations of the type $X \bowtie L = R$ and $L \bowtie Y = R$, where L, R are given known languages, X, Y are unknown languages, and \bowtie is the word blending operation.

Proposition 18. *The existence of a solution Y to the equation $L \bowtie Y = R$ is decidable for given regular languages L and R .*

Proof. According to [22], since \bowtie^r is the right-inverse of word blending, if there exists a solution Y to the given equation, then $Y' = (L \bowtie^r R^c)^c$ is also a solution. Moreover, in this case Y' is the maximal solution, in the sense that it includes all the other solutions to the equation. Since the family of regular languages is closed under \bowtie^r and complement, the algorithm for deciding the existence of a solution starts with constructing $L \bowtie Y'$, which is also regular, and checking whether $L \bowtie Y'$ equals R . As equality of regular languages is decidable [25], if the answer to the question “Is $L \bowtie Y'$ equal to R ?” is “yes”, then a solution to the equation exists, and Y' is such a solution. If the answer is “no”, then the equation has no solution. \square

Proposition 19. *The existence of a solution X to the equation $X \bowtie L = R$ is decidable for regular languages L and R .*

Proposition 20. *The existence of a singleton solution $\{w\}$ to the equation $L \bowtie \{w\} = R$ is decidable for regular languages L and R .*

Proof. If R is empty, a singleton solution $\{w\}$ to the equation $L \bowtie \{w\} = R$ exists if and only if L does not use all the letters from the alphabet Σ . The decision algorithm will check the emptiness of all regular languages $L \cap \Sigma^* a \Sigma^*$, where $a \in \Sigma$: If any of them is empty, then $\{w\} = \{a\}$ is a singleton solution, otherwise no singleton solution exists.

We now consider the case when R is not empty. If there is a singleton solution $\{w\}$ to the equation $L \bowtie \{w\} = R$, where $L, R \subseteq \Sigma^+$, $w \in \Sigma^+$ then there is a shortest singleton solution of length $k \geq 1$, denoted by $w_s = a_1 a_2 \cdots a_k$, with $a_1, a_2, \dots, a_k \in \Sigma$. We now want to show that the number of states in any finite state automaton that accepts R is at least k .

If $lg(w_s) = 1$, then $\lambda \notin R$, so the number of states of any finite state machine that recognizes R is at least 2, which is greater than the length of w_s .

Suppose $k \geq 2$. Define $L_i = (L \bowtie a_i) a_{i+1} \cdots a_k$ for $1 \leq i < k$, and define $L_k = L \bowtie a_k$. Then, we have $R = \bigcup_{i=1}^k L_i$. Note that $L_1 \not\subseteq \bigcup_{i=2}^k L_i$, as otherwise $a_2 a_3 \cdots a_k$ would be a shorter singleton solution than w_s —a contradiction.

Let $\alpha \in L_1 \subseteq R$; α can be represented as $\alpha = \alpha_1 a_1 a_2 \cdots a_k$, where $\alpha_1 \in \Sigma^*$. Assume now that R is recognized by a DFA $M = (Q, \Sigma, \delta, q_0, F)$ with $n < k$ states. Then there is a derivation

$$q_0 \alpha_1 a_1 a_2 \cdots a_k \Longrightarrow^* q_{i_1} a_1 a_2 \cdots a_k \Longrightarrow q_{i_2} a_2 \cdots a_k \Longrightarrow \cdots \Longrightarrow q_{i_k} a_k \Longrightarrow q_{i_{k+1}}.$$

Because M has $n < k$ states, there is a state that occurs twice in the set $\{q_{i_2}, q_{i_3}, \dots, q_{i_{k+1}}\}$.

If $q_{i_j} = q_{i_{k+1}}$ where $2 \leq j \leq k$, then $\alpha_1 a_1 \cdots a_{j-1} (a_j \cdots a_k)^+ \subseteq R$, and so there exists a word $\alpha_2 \in \Sigma^*$ such that $\alpha_1 a_1 \cdots a_{j-1} (a_j \cdots a_k)^+ \alpha_2 \subseteq L$. Thus, we have $\alpha \in \alpha_1 a_1 \cdots a_{j-1} (a_j \cdots a_k)^+ \alpha_2 \bowtie a_k \subseteq L_k \subseteq \bigcup_{i=2}^k L_i$.

If $q_{i_j} = q_{i_h}$ where $2 \leq j < h \leq k$, then $\alpha_1 a_1 \cdots a_{j-1} (a_j \cdots a_{h-1})^+ a_h \cdots a_k \subseteq R$, and so there exists a word $\alpha_2 \in \Sigma^*$ such that $\alpha_1 a_1 \cdots a_{j-1} (a_j \cdots a_{h-1})^+ \alpha_2 \subseteq L$. Then $\alpha \in (\alpha_1 a_1 \cdots a_{j-1} (a_j \cdots a_{h-1})^+ \alpha_2 \bowtie a_{h-1}) a_h \cdots a_k \subseteq L_{h-1} \subseteq \bigcup_{i=2}^k L_i$.

In either case, for all words $\alpha \in L_1$, $\alpha \in \bigcup_{i=2}^k L_i$. Thus, we have that $L_1 \subseteq \bigcup_{i=2}^k L_i$, which is a contradiction.

For the equation $L \bowtie Y = R$, if there is a singleton solution, there is a singleton solution w_s of minimal length k , and the number of states in any finite state machine for R is at least k . If the minimal deterministic finite automaton that generates R has k states, the algorithm for deciding the existence of a singleton solution will check all the words β , where $lg(\beta) \leq k$. The answer is “yes” if this algorithm finds a string β such that $L \bowtie \{\beta\} = R$, and “no” otherwise. \square

Proposition 21. *The existence of a singleton solution $\{w\}$ to the equation $\{w\} \bowtie L = R$ is decidable for regular languages L and R .*

Proposition 22. *The existence of a singleton solution $\{w\}$ to the equation $L \bowtie \{w\} = R$ is undecidable for regular languages R and context-free languages L .*

Proof. Assume, for the sake of contradiction, that the existence of a singleton solution $\{w\}$ to the equation $L \bowtie \{w\} = R$ is decidable for regular languages R and context-free languages L .

Given an arbitrary context-free language L' over an alphabet Σ , the context-free language $L_1 = \#\Sigma^+\# \cup L'\$$ can be constructed where $\#, \$ \notin \Sigma$. Note now that the equation $L_1 \bowtie \{w\} = \Sigma^*\$$ has a singleton solution $\{w\}$ if and only if $L' = \Sigma^*$ and the solution is $\{w\} = \{\$\}$. Thus, if we could decide the problem in the proposition, we would be able to decide whether or not $L' = \Sigma^*$ for arbitrary context-free languages L' , which is impossible. \square

Corollary 23. *The existence of a solution Y to the equation $L \bowtie Y = R$ is undecidable for regular languages R and context-free languages L .*

Proposition 24. 1. *The existence of a singleton solution $\{w\}$ to the equation $\{w\} \bowtie L = R$ is undecidable for a regular language R and a context-free language L .*

2. *The existence of a solution X to the equation $X \bowtie L = R$ is undecidable for a regular language R and a context-free language L .*

5 State Complexity

By Proposition 2, the family of regular languages is closed under word blending. Thus, we can consider the state complexity of the blending operation on two

regular languages. Recall from Proposition 2 that the blending of two languages can be expressed as a series of union, catenation, and quotient operations. While the state complexity of each of these operations is known, the state complexity of a combination of operations is not necessarily the same as the composition of the state complexities of the operations [29].

First, for illustrative purposes, we will construct an NFA that recognizes the blending of two languages given by DFAs. Let $A_m = (Q_m, \Sigma, \delta_m, s_m, F_m)$ be a DFA with $m \geq 1$ states that recognizes the language L_m and let $A_n = (Q_n, \Sigma, \delta_n, s_n, F_n)$ be a DFA with $n \geq 1$ states that recognizes the language L_n . We construct an NFA $B' = (Q', \Sigma, \delta', s', F')$, where $Q' = Q_m \cup Q_n$, $s' = s_m$, $F' = F_n$, and the transition function $\delta' : Q' \times \Sigma \rightarrow 2^{Q'}$ is defined for all $q \in Q'$ and $a \in \Sigma$ by

$$\delta'(q, a) = \begin{cases} \bigcup_{p \in Q_n} \delta_n(p, a) & \text{if } q \in Q_m \text{ and } \delta_m(q, a) \text{ is not the sink state,} \\ \delta_m(q, a) & \text{if } q \in Q_m \text{ and } \delta_m(q, a) \text{ is the sink state,} \\ \delta_n(q, a) & \text{if } q \in Q_n. \end{cases}$$

In Figure 1, we define two DFAs A_m and A_n and show the NFA B' resulting from the construction described above. Intuitively, the machine B' operates by first reading the input word assuming that it is the prefix of some word recognized by A_m . Since the blending occurs on only one symbol, the machine guesses at which symbol the blend occurs. Once the blend occurs the machine continues and assumes the rest of the word is the suffix of some word recognized by A_n .

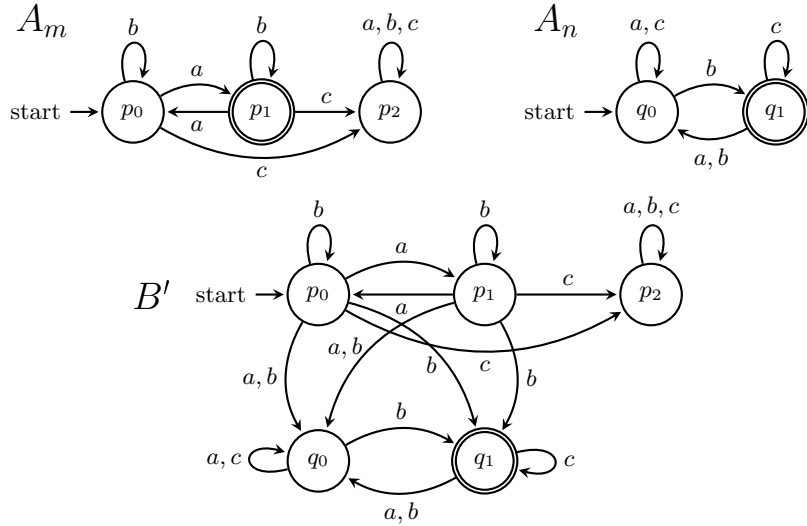


Fig. 1. The NFA B' recognizes the blend of the languages recognized by the DFAs A_m and A_n

Proposition 25. *The NFA B' recognizes the language $L_m \bowtie L_n$.*

Now, using the same basic idea, we will construct a DFA that recognizes the language of the blending of the two languages recognized by two given DFAs A_m and A_n . We construct a DFA $A' = (Q', \Sigma, \delta', s', F')$ where

- $Q' = Q_m \times 2^{Q_n}$,
- $s' = (s_m, \emptyset)$,
- $F' = \{(q, P) \in Q_m \times 2^{Q_n} \mid P \cap F_n \neq \emptyset\}$,
- $\delta'((q, P), a) = (\delta_m(q, a), P')$ for $a \in \Sigma$, where

$$P' = \begin{cases} \bigcup_{p \in P} \delta_n(p, a) & \text{if } \delta_m(q, a) \text{ is the sink state,} \\ \bigcup_{p \in Q_n} \delta_n(p, a) & \text{otherwise.} \end{cases}$$

Figure 2 shows the DFA A' that results from following the construction described above, where A_m and A_n are the DFAs shown in Figure 1. Each state of A' is a pair consisting of a state of A_m and a subset of states of A_n . Informally, we can divide the computation of a word into two phases. In the first phase, states of the form (q, P) are reached where q is not the sink state of A_m . Here, the set P is determined solely by the input symbol as the machine tries to guess the symbol on which the blending occurs. In the second phase, the machine reaches states (q_\emptyset, P) , where q_\emptyset is the sink state of A_m . The second phase only occurs when the blend occurs and the input that has been read is no longer a prefix of a word recognized by A_m . In this phase, the set P is determined by the transition function of A_n . We will show this formally in the following.

Proposition 26. *The DFA A' recognizes the language $L_m \bowtie L_n$.*

A simple count of the number of states in the state set of A' gives us as many as $m2^n$ states. We will show that, depending on the size of the alphabet, not all of these states are necessarily reachable. First, we consider the case where the alphabet is unary.

Theorem 27. *Let L_m and L_n be regular languages defined over a unary alphabet such that L_m is recognized by an m -state DFA and L_n is recognized by an n -state DFA. Then the state complexity of $L_m \bowtie L_n$ is $m + n - 1$ if both L_m and L_n are finite or 1 otherwise. Furthermore, this bound is reachable.*

Now, we will consider the state complexity when the languages are defined over alphabets of size greater than 1.

Lemma 28. *The DFA A' requires at most $(m - 1) \cdot (k - 1) + 2^n + 1$ states, where $k = |\Sigma| \leq 2^n$.*

Lemma 29. *Let $k \geq 3$ and $m, n \geq 2$. There exist families of DFAs A_m with m states and B_n with n states defined over an alphabet with k letters such that a DFA recognizing $A_m \bowtie B_n$ requires at least $(m - 1) \cdot (k - 1) + 2^n + 1$ states.*

These results together give us the following theorem.

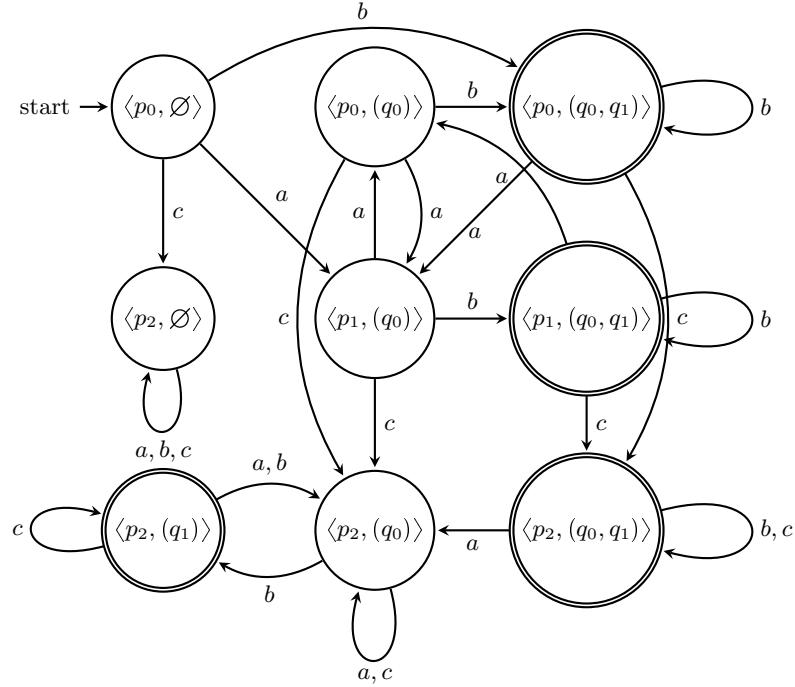
A' 

Fig. 2. The DFA A' recognizes the blend of the languages recognized by A_m and A_n from Figure 1

Theorem 30. Let A_m be a DFA with m states recognizing the language L_m and let A_n be a DFA with n states recognizing the language L_n , where L_m and L_n are defined over an alphabet Σ of size k . Then

$$sc(L_m \bowtie L_n) \leq (m - 1) \cdot (k - 1) + 2^n + 1,$$

and this bound can be reached in the worst case.

Acknowledgements. We thank Giuditta Franco for fruitful discussions on modelling the outcomes of various XPCR experiments.

References

1. Bonizzoni, P., Ferretti, C., Mauri, G., Zizza, R.: Separating some splicing models. *Information Processing Letters* 79(6), 255 – 259 (2001)
2. Bonizzoni, P., Felice, C.D., Zizza, R.: The structure of reflexive regular splicing languages via Schützenberger constants. *Theoretical Computer Science* 334(1), 71–98 (2005)

3. Brzozowski, J., Kari, L., Li, B., Szykuła, M.: State complexity of overlap assembly. arXiv preprint arXiv:1710.06000 (2017)
4. Carausu, A., Paun, G.: String intersection and short concatenation. *Revue Roumaine de Mathématiques Pures et Appliquées* 26(5), 713–726 (1981)
5. Csuhaaj-Varju, E., Petre, I., Vaszil, G.: Self-assembly of strings and languages. *Theoretical Computer Science* 374(1), 74–81 (2007)
6. Domaratzki, M.: Minimality in template-guided recombination. *Information and Computation* 207(11), 1209–1220 (2009)
7. Enaganti, S.K., Ibarra, O.H., Kari, L., Kopecki, S.: On the overlap assembly of strings and languages. *Natural Computing* 16(1), 175–185 (2017)
8. Enaganti, S.K., Ibarra, O.H., Kari, L., Kopecki, S.: Further remarks on DNA overlap assembly. *Information and Computation* 253, 143–154 (2017)
9. Franco, G.: A polymerase based algorithm for SAT. In: Coppo, M., Lodi, E., Pinna, G.M. (eds.) *Proc. Theoretical Computer Science, (ICTCS 2005)*. LNCS, vol. 3701, pp. 237–250. Springer, Berlin, Heidelberg (2005)
10. Franco, G., Bellamoli, F., Lampis, S.: Experimental analysis of XPCR-based protocols. arXiv preprint arXiv:1712.05182 (2017)
11. Franco, G., Giagulli, C., Laudanna, C., Manca, V.: DNA extraction by XPCR. In: Ferretti, C., Mauri, G., Zandron, C. (eds.) *Proc. DNA Computing, (DNA 10)*. LNCS, vol. 3384, pp. 104–112 (2005)
12. Franco, G., Manca, V.: Algorithmic applications of XPCR. *Natural Computing* 10(2), 805–819 (2011)
13. Franco, G., Manca, V., Giagulli, C., Laudanna, C.: DNA recombination by XPCR. In: Carbone, A., Pierce, N.A. (eds.) *Proc. DNA Computing, (DNA 11)*. LNCS, vol. 3892, pp. 55–66 (2006)
14. Golan, J.S.: *The Theory of Semirings with Applications in Mathematics and Theoretical Computer Science*. Addison-Wesley Longman Ltd. (1992)
15. Goode, E., Pixton, D.: Recognizing splicing languages: Syntactic monoids and simultaneous pumping. *Discrete Applied Mathematics* 155(8), 989–1006 (2007)
16. Gries, S.T.: Shouldn't it be breakfunch? A quantitative analysis of blend structure in English. *Linguistics* 42(3), 639–667 (2004)
17. Head, T.: Formal language theory and DNA: An analysis of the generative capacity of specific recombinant behaviors. *Bulletin of Math. Biology* 49(6), 737–759 (1987)
18. Holzer, M., Jakobi, S.: Chop operations and expressions: Descriptive complexity considerations. In: Mauri, G., Leporati, A. (eds.) *Developments in Language Theory*. LNCS, vol. 6795, pp. 264–275. Springer, Berlin, Heidelberg (2011)
19. Holzer, M., Jakobi, S.: State complexity of chop operations on unary and finite languages. In: Kutrib, M., Moreira, N., Reis, R. (eds.) *Descriptive Complexity of Formal Systems*. LNCS, vol. 7386, pp. 169–182. Springer, Berlin, Heidelberg (2012)
20. Holzer, M., Jakobi, S., Kutrib, M.: The chop of languages. In: *Automata and Formal Languages, 13th International Conference, AFL 2011, Debrecen*. pp. 197–210 (2011)
21. Ito, M., Lischke, G.: Generalized periodicity and primitivity for words. *Mathematical Logic Quarterly* 53(1), 91–106 (2007)
22. Kari, L.: On language equations with invertible operations. *Theoretical Computer Science* 132(1-2), 129–150 (1994)
23. Manca, V., Franco, G.: Computing by polymerase chain reaction. *Mathematical Biosciences* 211(2), 282–298 (2008)
24. Mateescu, A., Păun, G., Rozenberg, G., Salomaa, A.: Simple splicing systems. *Discrete Applied Mathematics* 84(1-3), 145–163 (1998)
25. Mateescu, A., Salomaa, A.: *Handbook of Formal Languages*. Springer-Verlag New York, Inc., New York, NY (1997)

26. Pixton, D.: Splicing in abstract families of languages. *Theoretical Computer Science* 234, 135–166 (2000)
27. Păun, G.: On the splicing operation. *Discrete Applied Mathematics* 70(1), 57–79 (1996)
28. Salomaa, A.: *Formal Languages*. Academic Press, Inc., New York, NY (1977)
29. Salomaa, K., Yu, S.: On the state complexity of combined operations and their estimation. *International Journal of Foundations of Computer Science* 18(4), 683–698 (2007)

A Appendix

Here we include proofs that were omitted in the paper due to the limitation on the number of pages.

Lemma 1. *If x, y are non-empty words in Σ^+ , then* (p. 4)

$$x \bowtie y = \{z \in \Sigma^+ \mid \exists \alpha, \beta, \gamma_1, \gamma_2 \in \Sigma^*, \exists a \in \Sigma : x = \alpha a \gamma_1, y = \gamma_2 a \beta, z = \alpha a \beta\}.$$

Proof. Let A denote the right hand side of the equality. The inclusion $A \subseteq x \bowtie y$ is obvious by the definition of word blending. To prove the converse, let $z \in x \bowtie y$. Then $z = \alpha w \beta$ for some $\alpha, \beta, \gamma_1, \gamma_2 \in \Sigma^*$ and $w \in \Sigma^+$ such that $x = \alpha w \gamma_1$, $y = \gamma_2 w \beta$. As $w \in \Sigma^+$, w can be written as $w = w_1 a$, where $w_1 \in \Sigma^*$, $a \in \Sigma$. It follows that $x = \alpha w_1 a \gamma_1$, $y = \gamma_2 w_1 a \beta$ and $z = \alpha w_1 a \beta$, that is, $x = \alpha' a \gamma_1$, $y = \gamma_2' a \beta$ and $z = \alpha' a \beta$ with $\alpha' = \alpha w_1 \in \Sigma^*$ and $\gamma_2' = \gamma_2 w_1 \in \Sigma^*$. Thus, $z \in A$, and the equality is proved. \square

Proposition 2. *Given languages $L_1, L_2 \subseteq \Sigma^+$,* (p. 4)

$$L_1 \bowtie L_2 = \bigcup_{a \in \Sigma} (L_1 (a \Sigma^*)^{-1}) a ((\Sigma^* a)^{-1} L_2)$$

Proof. (\subseteq) Let $z \in L_1 \bowtie L_2$. Then, by Lemma 1, $z = \alpha a \beta$, for some $x \in L_1$ and $y \in L_2$ such that $x = \alpha a \gamma_1$, $y = \gamma_2 a \beta$ where $a \in \Sigma, \alpha, \beta, \gamma_1, \gamma_2 \in \Sigma^*$. It is clear that $\alpha \in L_1 (a \Sigma^*)^{-1}$ and $\beta \in (\Sigma^* a)^{-1} L_2$, so $z = \alpha a \beta \in (L_1 (a \Sigma^*)^{-1}) a ((\Sigma^* a)^{-1} L_2)$.

(\supseteq) Let $z \in \bigcup_{a \in \Sigma} (L_1 (a \Sigma^*)^{-1}) a ((\Sigma^* a)^{-1} L_2)$. Then there exists $a \in \Sigma$ and words $\alpha, \gamma_1, \gamma_2, \beta \in \Sigma^*$, such that $z = \alpha a \beta$, where $x = \alpha a \gamma_1 \in L_1, y = \gamma_2 a \beta \in L_2$, which implies that $z \in L_1 \bowtie L_2$. \square

Corollary 3. *Every full AFL is closed under word blending.* (p. 4)

Proof. It follows from Proposition 2, and every full AFL is closed under left/right quotient with regular languages, catenation and finite union [28]. \square

Proposition 8. *The family of context-sensitive languages is not closed under the right inverse of blending.* (p. 5)

Proof. Let L_0 be a recursively enumerable language over Σ , that is not context-sensitive, and the context-sensitive language L_1 over $\Sigma \cup \{a, b\}$ with $a, b \notin \Sigma$, that can be associated to L_0 such that L_1 consists of words of the form $a^i b P$ where $i \geq 0$ and $P \in L_0$ and, in addition, for every $P \in L_0$ there is an $i \geq 0$ such that $a^i b P \in L_1$.

The result now follows as $(L_1 \bowtie^r \{a^* b\}) \cap \Sigma^* b \Sigma^* = \{\Sigma^* b P \mid b \notin \Sigma, P \in L_0\}$, which is not context-sensitive, and the family of context-sensitive languages is closed under intersection with regular languages. \square

Proposition 9. *The left inverse of blending can be expressed using the right inverse of blending, and mirror image as $w \bowtie^l v = \text{mi}(\text{mi}(v) \bowtie^r \text{mi}(w))$.* (p. 6)

Proof. If $w \in x \bowtie v$, there exist $\alpha, \beta, \gamma_1, \gamma_2 \in \Sigma^*, b \in \Sigma$ such that $w = \alpha b \beta, x = \alpha b \gamma_1, v = \gamma_2 b \beta$ by Lemma 1. Then, since $x = \alpha b \gamma_1 = \text{mi}(\text{mi}(\gamma_1) b \text{mi}(\alpha))$, we have

$$\begin{aligned}
x &\in \text{mi}(\Sigma^* b \text{mi}(\alpha)) \\
&= \text{mi} \left(\Sigma^* b \left((\text{mi}(\beta) b)^{-1_l} (\text{mi}(\beta) b \text{mi}(\alpha)) \right) \right) \\
&= \text{mi} \left(\Sigma^* b \left(\left((\text{mi}(\beta) b \text{mi}(\gamma_2)) (b \text{mi}(\gamma_2))^{-1} b \right)^{-1_l} \text{mi}(w) \right) \right) \\
&= \text{mi} \left(\Sigma^* b \left(\left((\text{mi}(v) (b \text{mi}(\gamma_2))^{-1} b \right)^{-1_l} \text{mi}(w) \right) \right) \\
&\subseteq \text{mi} \left(\bigcup_{a \in \Sigma} \left(\Sigma^* a \left(\left((\text{mi}(v) (a \Sigma^*)^{-1} a \right)^{-1_l} \text{mi}(w) \right) \right) \right) \\
&= \text{mi}(\text{mi}(v) \bowtie^r \text{mi}(w)) \\
&= w \bowtie^l v.
\end{aligned}$$

Now consider $x \in w \bowtie^l v$. Then,

$$\begin{aligned}
x \in w \bowtie^l v &= \text{mi}(\text{mi}(v) \bowtie^r \text{mi}(w)) \\
&= \text{mi} \left(\bigcup_{a \in \Sigma} \Sigma^* a \left(\left((\text{mi}(v) (a \Sigma^*)^{-1} a \right)^{-1_l} \text{mi}(w) \right) \right) \right).
\end{aligned}$$

Thus, there exist $b \in \Sigma, \gamma_1 \in \Sigma^*, \gamma_3 \in (\text{mi}(v) (b \Sigma^*)^{-1} b)$ such that

$$x = \text{mi}(\text{mi}(\gamma_1) b (\gamma_3^{-1_l} \text{mi}(w))) = (w \text{mi}(\gamma_3)^{-1} b) \gamma_1.$$

Then we have

$$\begin{aligned}
w &\in (w \text{mi}(\gamma_3)^{-1} b) ((\Sigma^* b)^{-1_l} v) \\
&\subseteq (w \text{mi}(\gamma_3)^{-1} b) \gamma_1 (b \Sigma^*)^{-1} b ((\Sigma^* b)^{-1_l} v) \\
&\subseteq \bigcup_{a \in \Sigma} (((w \text{mi}(\gamma_3)^{-1} b) \gamma_1) (a \Sigma^*)^{-1} a) ((\Sigma^* a)^{-1_l} v) \\
&= \bigcup_{a \in \Sigma} (x (a \Sigma^*)^{-1} a) ((\Sigma^* a)^{-1_l} v) \\
&= x \bowtie v.
\end{aligned}$$

□

(p. 7)

Lemma 13. Let $L \subseteq \Sigma^+$ be a language. Then for each word $w \in L^{\bowtie^*}$, there exists $n \in \mathbb{N}$ such that there are words $u_i \in \text{inf}(L), 1 \leq i \leq n$ and $\alpha_j \in \Sigma^*, 1 \leq j \leq n$ and symbols $a_k \in \Sigma, 1 \leq k \leq n-1$ where

1. for $n > 1$,

- (a) $w = \alpha_1 a_1 \alpha_2 a_2 \cdots a_{n-1} \alpha_n$,
 (b) $u_i = a_{i-1} \alpha_i a_i \in \text{inf}(L)$ for all $2 \leq i \leq n-1$,
 (c) $u_1 = \alpha_1 a_1 \in \text{pref}(L)$ and $u_n = a_{n-1} \alpha_n \in \text{suff}(L)$,
2. $u_1 = w \in L$ for $n = 1$.

Proof. Let $w \in L^{\boxtimes*}$. Then $w \in L^{\boxtimes j}$ for some $j \geq 0$. We will prove the statement by induction on j . If $j = 0$ then $w \in L$ and the statement holds taking $n = 1$. Assume that the statement holds for words in $L^{\boxtimes j}$ for any $j \leq k$, and consider a word $w \in L^{\boxtimes k+1} = L \boxtimes L^{\boxtimes k}$. This implies that $w \in x \boxtimes y$ for some $x \in L$ and $y \in L^{\boxtimes k}$. By the induction hypothesis, either $y \in L$ or $y = \beta_1 b_1 \beta_2 b_2 \cdots b_{m-1} \beta_m$ for some $m \geq 2$ with $b_i \in \Sigma$, $1 \leq i \leq m-1$ and $\beta_j \in \Sigma^*$, $1 \leq j \leq m$ satisfying the conditions of the Lemma.

If $y \in L$, then $x \boxtimes y$ consists of all words of the form $\alpha_1 a_1 \alpha_2$ where $x = \alpha_1 a_1 \gamma_1$ for $\alpha_1, \gamma_1 \in \Sigma^*$ and $a_1 \in \Sigma$ and $y = \gamma_2 a_1 \alpha_2$ for some $\gamma_2, \alpha_2 \in \Sigma^*$. It is easy to see that $\alpha_1 a_1 \alpha_2$ satisfies the conditions of the Lemma.

Otherwise, if $y = \beta_1 b_1 \beta_2 b_2 \cdots b_{m-1} \beta_m$ for some $m \geq 2$, then the set $x \boxtimes y$ consists of words of the form $\alpha'_1 a'_1 \beta'_\ell b_\ell \cdots b_{m-1} \beta_m$ where $\alpha'_1 a'_1 \in \text{pref}(x)$ and $1 \leq \ell \leq m$. Here, we observe that in order for the blend to occur, we have $\alpha'_1 \beta'_\ell \in \text{suff}(b_{\ell-1} \beta_\ell)$. Then by definition, we have $\alpha'_1 a'_1 \in \text{pref}(x) \subseteq \text{pref}(L)$ and $\alpha'_1 \beta'_\ell b_\ell \in \text{inf}(L)$ and the rest follows. \square

Proposition 14. *Given an NFA A , there exists an NFA A' recognizing the language $L(A)^{\boxtimes*}$ which is effectively constructible.* (p. 7)

Proof. Given an NFA $A = (Q, \Sigma, \delta, s, F)$, we can construct an NFA $A' = (Q', \Sigma, \delta', s', F')$ that recognizes the language $L(A)^{\boxtimes*}$. Informally, the machine operates by guessing when a blend occurs. Recall from Lemma 13, words of $L(A)^{\boxtimes*}$ are of the form $\alpha_1 a_1 \alpha_2 a_2 \cdots a_{n-1} \alpha_n$. When a blend occurs on a symbol a_i , the machine then simulates the operation of A on the blended suffix α_{i+1} and continues to guess when the next blend may occur. This process repeats until the machine reaches a final state of A and accepts or it does not and the machine rejects.

Formally, we define A' by

- $Q' = \{\langle p \rangle, \langle q, r \rangle \mid p, q, r \in Q\}$,
- $s' = \langle s \rangle$,
- $F' = \{\langle q \rangle \mid q \in F\}$,

and the transition function is defined by

- $\delta'(\langle q \rangle, a) = \{\langle q' \rangle, \langle q', r' \rangle \mid q' \in \delta(q, a), r' \in \bigcup_{p \in Q} \delta(p, a)\}$,
- $\delta'(\langle q, r \rangle, a) = \{\langle r' \rangle, \langle r', p' \rangle \mid r' \in \delta(r, a), p' \in \bigcup_{p \in Q} \delta(p, a)\}$.

First, we show that $L(A') \subseteq L(A)^{\boxtimes*}$. Consider a word $w \in L(A')$. There exists a sequence of states, or path, in A' on w from $\langle s \rangle$ to $\langle q'_n \rangle$ where $\langle q'_n \rangle \in F'$. Recall that states of A' are of the form $\langle q \rangle$ or $\langle q, r \rangle$. Of the states on the path defined by the computation of w , we consider those states of the form $\langle q, r \rangle$ and

label them $\langle q'_i, q_{i+1} \rangle$ for $0 \leq i \leq n$. Each state $\langle q'_i, q_{i+1} \rangle$ is entered upon reading a symbol which we will call $a_i \in \Sigma$.

Now for each $1 \leq i \leq n-1$, consider the path in A' between $\langle q'_{i-1}, q_i \rangle$ and $\langle q'_i, q_{i+1} \rangle$. Between these two states, each state on the path is of the form $\langle q, r \rangle$, as we have already labeled states that are of the form $\langle q, r \rangle$. This implies that there is a word $\alpha_{i+1}a_{i+1}$ which takes A from q_i to q'_i . But this means that $a_i\alpha_{i+1}a_{i+1}$ is a subword of some word recognized by A . Then we can write $w = \alpha_1a_1\alpha_2a_2 \cdots a_{n-1}\alpha_n$ where $\alpha_1a_1 \in \text{pref}(L(A))$ since it takes A from s to s' and $\alpha_n \in \text{suff}(L(A))$ since it takes A from a state q_n to $q'_n \in F$. Thus, we have $w \in L(A)^{\bowtie*}$.

Now, we show that $L(A)^{\bowtie*} \subseteq L(A')$. Consider a word $w \in L(A)^{\bowtie*}$. We can write $w = \alpha_1a_1\alpha_2a_2 \cdots a_{n-1}\alpha_n$. Since each $a_i\alpha_{i+1}a_{i+1}$ is a subword of a word recognized by A , there must exist a path between two states of A , say q_i and q'_i , on each word $\alpha_{i+1}a_{i+1}$. Then a path can be traced in A' from the initial state $\langle s \rangle$ by

$$\langle s \rangle \xrightarrow{\alpha_1a_1} \langle q'_0, q_1 \rangle \xrightarrow{\alpha_2a_2} \langle q'_1, q_2 \rangle \xrightarrow{\alpha_3a_3} \cdots \xrightarrow{\alpha_{n-1}a_{n-1}} \langle q'_{n-1}, q_n \rangle \xrightarrow{\alpha_n} \langle q'_n \rangle.$$

Recall that by Lemma 13, $\alpha_1a_1 \in \text{pref}(L(A))$ and therefore there is a path from s to a state q'_0 in A . Note also that since $a_{n-1}\alpha_n \in \text{suff}(L(A))$, the state q'_n which is reached on a path on α_n must be an accepting state of A and therefore $\langle q'_n \rangle \in F'$ and $w \in L(A')$. \square

(p. 7) **Proposition 15.** *Let L be a regular language. It is decidable whether or not L is closed under \bowtie_* .*

Proof. Let A be an NFA that recognizes L . Then by the construction given in Proposition 14, we can construct an NFA A' that recognizes $L^{\bowtie*}$. Testing equivalence of two NFAs is known to be decidable [28] and therefore, testing whether $L = L^{\bowtie*}$ is decidable. \square

(p. 7) **Theorem 16.** *It is decidable whether or not a regular language has a base over \bowtie_* .*

Proof. Let $L \subseteq \Sigma^*$ be a regular language given as a finite automaton. Let $R = \{w \in \Sigma^* \mid |w|_a \leq 2 \text{ for all } a \in \Sigma\}$ be the set of words in which each symbol of Σ appears at most twice. We claim that if L is closed under \bowtie_* , it must be generated by a base $B \subseteq L \cap R$.

Suppose otherwise and that L is generated by a base B' which is not a subset of $L \cap R$. Then B' contains a word of the form $w = x_1ax_2ax_3ax_4$, where $x_1, x_2, x_3, x_4 \in \Sigma^*$ and $a \in \Sigma$. Let $w_1 = x_1ax_2ax_4$ and $w_2 = x_1ax_3ax_4$ and note that $w_1, w_2 \in w \bowtie w$ and therefore $w_1, w_2 \in L$. Furthermore, we have $w \in w_1 \bowtie w_2$ and we can define an equivalent base $B'' = (B' \setminus \{w\}) \cup \{w_1, w_2\}$.

Now, we show that we only need to repeat this procedure a finite number of times. One only needs to consider words of length at most n , where n is the pumping length of L . Indeed, consider a word u of length greater than n . Since L is regular, we can write $u = xy^2z$ where $x, y, z \in \Sigma^*$ and $xyz \in L$ is of length

at most n . But then we have $xyz \in u \bowtie u$. Thus, it suffices to consider only those words in L of length up to n .

We can test whether the base B obtained via this process generates L . Using the construction of Proposition 14, we can construct an NFA C that recognizes $B^{\bowtie*}$ and test whether $L(C) = L(A)$. This is decidable since NFA equivalence is known to be decidable [28]. \square

Proposition 19. *The existence of a solution X to the equation $X \bowtie L = R$ is decidable for regular languages L and R .* (p. 8)

Proof. Similar to the preceding proof, and using the closure of the family of regular languages under the left-inverse of word blending. \square

Proposition 21. *The existence of a singleton solution $\{w\}$ to the equation $\{w\} \bowtie L = R$ is decidable for regular languages L and R .* (p. 9)

Proof. Similar to the preceding proof, and using the fact that the minimal length of a singleton solution to the equation is equal to or smaller than the number of states of any deterministic finite automaton that recognizes R . \square

Proposition 25. *The NFA B' recognizes the language $L_m \bowtie L_n$.* (p. 11)

Proof. First, we show that $L_m \bowtie L_n \subseteq L(B')$. Let $w \in L_m \bowtie L_n$ and write $w = \alpha a \beta$ for a symbol $a \in \Sigma$ and words $\alpha, \beta \in \Sigma^*$ such that for some words $\gamma_1, \gamma_2 \in \Sigma^*$, we have $\alpha a \gamma_1 \in L_m$ and $\gamma_2 a \beta \in L_n$. Since αa is a prefix of a word in L_m , let $p = \delta_m(s_m, \alpha a) \in \delta'(s', \alpha a)$. However, since $a \beta$ is a suffix of a word of L_n , there exists a state $q \in Q_n$ such that $\delta_n(q, a) = r$ and by the definition of δ' , we have $r \in \delta'(s', \alpha a)$. From here, we observe that we must have $\delta_n(r, \beta) \in F_n$ and therefore $\delta'(r, \beta) \in F'$ and w is accepted by B' .

Next, we show that $L(B') \subseteq L_m \bowtie L_n$ and consider a word $w \in L(B')$. By definition, must exist a path on w from $s' = s_m$ to a state in F_n and we can divide the path into two parts. The first part consists of transitions among states of A_m and the latter part consists of transitions among states of Q_n . Observe that the only way for a transition from a state $p \in Q_m$ to a state $q \in Q_n$ to be defined is if for some symbol $a \in \Sigma$, $\delta_m(p, a)$ is not a transition to a sink state and that $\delta_n(r, a) = q$ for some $r \in Q_n$. Thus, we can write $w = xay$ for words $x, y \in \Sigma^*$, where $a \in \Sigma$ is the symbol on which the path transitions from states of A_m to states of A_n . Then this implies that xa is a prefix of some word in L_m and ay is a suffix of some word in L_n . Therefore, $w \in L_m \bowtie L_n$ by definition and thus, we have shown that B' recognizes $L_m \bowtie L_n$. \square

Proposition 26. *The DFA A' recognizes the language $L_m \bowtie L_n$.* (p. 11)

Proof. First, to show that $L_m \bowtie L_n \subseteq L(A')$, consider a word $w \in L_m \bowtie L_n$. Then $w = \alpha a \beta$ for some symbol $a \in \Sigma$ and words $\alpha, \beta \in \Sigma^*$ where for some $\gamma_1, \gamma_2 \in \Sigma^*$, we have $\alpha a \gamma_1 \in L_m$ and $\gamma_2 a \beta \in L_n$. Observe that since αa is a prefix of a word in L_m , the state $\delta_m(s_m, \alpha a)$ is not the sink state. Similarly, since $a \beta$

is the suffix of a word in L_n , there exists at least one state in Q_n that has an incoming transition on the symbol a .

If $\beta = \lambda$, then $a \in \text{suff}(L_n)$ and $\delta'(s', \alpha a) \in F'$ and therefore, $w \in L(A')$. So suppose $\beta = \sigma\beta'$ for some symbol $\sigma \in \Sigma$ and $\beta' \in \Sigma^*$. We assume that $\delta_m(s_m, \alpha a \sigma)$ is the sink state, since otherwise, we can write $w = \alpha'\sigma\beta'$ where $\alpha' = \alpha a$ and repeat the same process. Then reading σ takes us from $(\delta_m(s_m, \alpha a), P)$ to the state (q_\emptyset, P') , where q_\emptyset denotes the sink state of A_m and $P' = \bigcup_{p \in P} \delta_n(p, a)$. Since β is the suffix of a word in L_n , there exists a state $p \in P$ such that $\delta_n(p, \beta) \in F_n$. Thus, reading the rest of β takes us to a state (q_\emptyset, P'') with $P'' \cap F_n \neq \emptyset$ and therefore, w is recognized by A' .

To show that $L(A') \subseteq L_m \bowtie L_n$, we consider a word w recognized by A' . That is, upon reading w , the machine A' reaches a final state (q', P') with $P' \cap F_n \neq \emptyset$. First, suppose that q' is not the sink state of A_m . We can write $w = w'a$ for a symbol $a \in \Sigma$ and a word $w' \in \Sigma^*$. Since q' is not the sink state of A_m , the word w is a prefix of some word in L_m and we have $w'a\gamma_1 \in L_m$ for some $\gamma_1 \in \Sigma^*$. By the definition of the transition function, a is a suffix of a word in L_n and we have $\gamma_2 a \in L_n$ for some $\gamma_2 \in \Sigma^*$. Thus, $w \in w'a\gamma_1 \bowtie \gamma_2 a$ and therefore $w \in L_m \bowtie L_n$.

Now, suppose that q' is the sink state of A_m . Let $w = \alpha ab\beta$ such that αab is the shortest prefix of w that enters the sink state q' of A_m . Then αa is a prefix of a word in L_m so we have $\alpha a\gamma_1 \in L_m$ for some $\gamma_1 \in \Sigma^*$. Reading αa takes us to the state $\delta'(s', \alpha a) = (q, P)$ such that q is some state of A_m which is not the sink state and $P = \bigcup_{p \in Q_n} \delta_n(p, a)$.

We claim that $ab\beta$ is a suffix of a word in L_n . To see this, we observe that since reading b from (q, P) takes A' to the state (q', P'') , where q' is the sink state of A_m , any transitions from P'' no longer depend solely on the input symbol. Therefore, there must exist a path in A_n from a state $r \in P$ to a final state of A_n on the word $ab\beta$. We can write $\beta' = b\beta$ and thus there exists a word γ_2 such that $\gamma_2 a\beta' \in L_n$. Therefore, $w \in \alpha a\gamma_1 \bowtie \gamma_2 a\beta'$ and $w \in L_m \bowtie L_n$ as desired.

Thus, we have shown that $L(A') = L_m \bowtie L_n$. \square

(p. 11)

Theorem 27. *Let L_m and L_n be regular languages defined over a unary alphabet such that L_m is recognized by an m -state DFA and L_n is recognized by an n -state DFA. Then the state complexity of $L_m \bowtie L_n$ is either $m + n - 1$ if both L_m and L_n are finite, or 1 otherwise. Furthermore, this bound is reachable.*

Proof. Recall that by Proposition 2, $L_m \bowtie L_n = (L_m(a^+)^{-1})a((a^+)^{-1}L_n)$. If either L_m or L_n are infinitely large, then we have $L_m \bowtie L_n = a^*$, in which case the state complexity of $L_m \bowtie L_n$ is 1. If both L_m and L_n are finite, then it is easy to see that the state complexity of $L_m \bowtie L_n$ is $m + n - 1$. \square

(p. 11)

Lemma 28. *The DFA A' requires at most $(m - 1) \cdot (k - 1) + 2^n + 1$ states, where $k = |\Sigma| \leq 2^n$.*

Proof. First, observe that in order to maximize the number of reachable states of A' , the DFA A_m must contain a state that cannot reach an accepting state. Otherwise, if every state of A_m can reach an accepting state, then by definition

of A' , we have $L(A_m) \bowtie L(A_n) \subseteq \text{pref}(L(A_m))$. One can construct a DFA for $\text{pref}(L(A_m))$ by modifying A_m so that every state of A_m is a final state. In this case, A' would then require at most m states. Thus, we assume that A_m contains a sink state q_\emptyset which cannot reach an accepting state.

Consider the transition function δ' on a state (q, P) , where $q \neq q_\emptyset$. Then for each symbol $a \in \Sigma$, there is only one possible reachable set of states P in A_n . This gives us up to $(m-1) \cdot k$ reachable states. However, we claim that in order for two states (q, P) and (q, P') with $P \neq P'$ to be distinguishable, q must contain a transition to q_\emptyset . Otherwise, for every symbol $a \in \Sigma$, we have $\delta'((q, P), a) = \delta'((q, P'), a)$ by definition. Thus, since every state must contain at least one transition to q_\emptyset and A_m is deterministic, A' has only at most $(m-1) \cdot (k-1)$ reachable states of this form.

Next, consider that there are up to 2^n reachable states (q_\emptyset, P) as derived from the subset construction.

Finally, we note that the initial state $s' = (s_m, \emptyset)$ does not belong to any of the above sets. Adding all of these states together, we have at most $(m-1) \cdot (k-1) + 2^n + 1$ reachable states. \square

Lemma 29. *Let $k \geq 3$ and $m, n \geq 2$. There exist families of DFAs A_m with m states and B_n with n states defined over an alphabet with k letters such that a DFA recognizing $A_m \bowtie B_n$ requires at least $(m-1) \cdot (k-1) + 2^n + 1$ states.* (p. 11)

Proof. Let $\Sigma = \{a_1, \dots, a_{k-2}, b, c\}$. We will define the DFAs A_m and B_n over Σ .

Let $A_m = (Q_m, \Sigma, \delta_m, s_m, F_m)$ where $Q_m = \{0, \dots, m-1\}$, $s_m = 0$, and $F_m = \{m-2\}$. We define the transition function δ_m by

- $\delta_m(p, a_i) = p$ for all $0 \leq p \leq m-2$ and $1 \leq i \leq k-2$,
- $\delta_m(p, b) = p+1$ for $0 \leq p \leq m-2$,
- $\delta_m(m-1, \sigma) = m-1$ for all $\sigma \in \Sigma$.

The DFA A_m is shown in Figure 3.

Let $B_n = (Q_n, \Sigma, \eta_n, s_n, F_n)$ where $Q_n = \{0, \dots, n-1\}$, $s_n = 0$, and $F_n = \{n-1\}$. We will define the transition function η_n by

- $\eta_n(q, b) = q+1 \pmod n$ for $0 \leq q \leq n-1$,
- $\eta_n(q, c) = q$ for $0 \leq q \leq n-1$.

For transitions on symbols a_i with $1 \leq i \leq k-2$, we define an enumeration of the subsets of Q_n and let $Q_n[i]$ be the i th subset of Q_n . Any arbitrary enumeration of subsets of Q_n suffices for this proof subject to the condition that

1. for $0 \leq i \leq k-2$, each i corresponds to a particular subset of Q_n and
2. we reserve $Q_n[0] = Q_n$ and $Q_n[1] = \{0, 1, \dots, n-2\}$.

That is, $Q_n[i] \neq Q_n[j]$ iff $i \neq j$ for $0 \leq i, j \leq k-2$. Also note that while we have defined $Q_n[0]$, there is no symbol a_0 . We will show later that, by our definitions, the role of a_0 will be played by b . If $k > 2^n$, then this property cannot hold but it is clear that we can enumerate all 2^n subsets of Q_n .

Then we define transitions on $a_i \in \Sigma$ by

$$\eta(q, a_i) = \begin{cases} q & \text{if } q \in Q_n[i], \\ q + \min_{(q+j \bmod n) \in Q_n[i]} j \bmod n & \text{if } q \notin Q_n[i]. \end{cases}$$

In other words, for each state $q \in Q_n$, the transition on the symbol a_i goes to the “next” state that is in $Q_n[i]$. If $q \in Q_n[i]$, then that q itself is the “next” state.

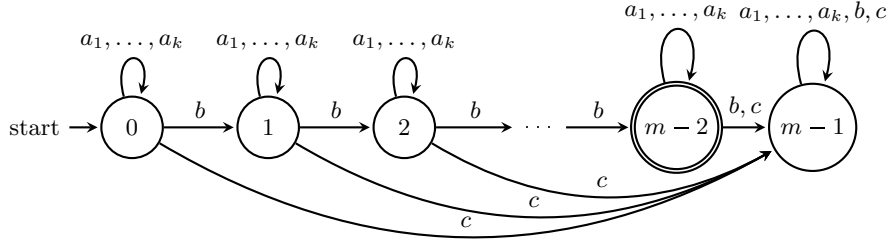


Fig. 3. The DFA A_m .

We will show that A' contains $(m-1) \cdot (k-1) + 2^n + 1$ reachable and distinguishable states.

First, to show that the states are reachable, we note that $s' = (s_m, \emptyset)$ is clearly reachable as the initial state. Then, we observe that for $1 \leq i \leq k-2$, the state $(q, Q_n[i])$ with $q \in Q_n \setminus \{m-1\}$ is reachable on the word $b^q a_i$ and $(q, Q_n[0])$ is reachable on the word b^q . Since the only symbol not used here is c , this gives us $(m-1) \cdot (k-1)$ states.

Now we consider states of the form $(m-1, P)$ where $P \subseteq Q_n$. Observe that $(m-1, Q_n)$ can be reached on the word b^{m-1} . Then, we note that for an arbitrary set $T \subseteq Q_n$ and element $t \in T$, we can reach the state $(m-1, T \setminus \{t\})$ from the state $(m-1, T)$ on the word $b^{n-t} a_1 b^t$. By repeating this process, we can reach any state (q, P) for any arbitrary subset $P \subseteq Q_n$. Thus, we have an additional 2^n reachable states of the form $(m-1, P)$, giving us a total of $(m-1) \cdot (k-1) + 2^n + 1$ reachable states.

Next, we will show that these states are pairwise distinguishable. Consider two states (q, P) and (q', P') . First, we fix $P = P'$ and assume without loss of generality that $q < q'$. Then the two states are distinguished by the word $b^{m-1-q} a_1^n$.

Now, we consider when $P \neq P'$. In this case, reading c takes the state (q, P) to $(m-1, P)$ and (q', P') to $(m-1, P')$. Then without loss of generality, there exists an element $t \in P$ and $t \notin P'$. Then these states are distinguished by the word b^{n-t} .

Thus, we have shown that all $(m-1) \cdot (k-1) + 2^n + 1$ states are reachable and pairwise distinguishable. \square