

# Closure Properties of Pattern Languages

Joel D. Day<sup>1</sup>, Daniel Reidenbach<sup>1</sup>, and Markus L. Schmid<sup>2</sup>

<sup>1</sup> Department of Computer Science, Loughborough University,  
Loughborough, Leicestershire, LE11 3TU, UK

J.Day@lboro.ac.uk,

D.Reidenbach@lboro.ac.uk

<sup>2</sup> Universität Trier, FB IV–Abteilung Informatikwissenschaften,  
D-54286 Trier, Germany

MSchmid@uni-trier.de

**Abstract.** Pattern languages are a well-established class of languages that is particularly popular in algorithmic learning theory, but very little is known about their closure properties. In the present paper we establish a large number of closure properties of the terminal-free pattern languages, and we characterise when the union of two terminal-free pattern languages is again a terminal-free pattern language. We demonstrate that the equivalent question for general pattern languages is characterised differently, and that it is linked to some of the most prominent open problems for pattern languages. We also provide fundamental insights into a well-known construction of E-pattern languages as unions of NE-pattern languages, and vice versa.

**Keywords:** Pattern languages, Closure properties.

## 1 Introduction

Pattern languages were introduced by Dana Angluin [1] in order to model the algorithmic inferrability of patterns that are common to a set of words. In this context, a pattern is a sequence of variables and terminal symbols, and its language is the set of all words that can be generated from the pattern by a substitution that replaces all variables in the pattern by words of terminal symbols. Hence, more formally, a substitution is a terminal-preserving morphism, i. e., a morphism that maps every terminal symbol to itself. For example, the pattern language of the pattern  $\alpha := x_1x_1\mathbf{a}x_2\mathbf{b}$ , where  $x_1, x_2$  are variables and  $\mathbf{a}, \mathbf{b}$  are terminal symbols, is the set of all words that have a square as a prefix, followed by an arbitrary suffix that begins with the letter  $\mathbf{a}$  and ends with the letter  $\mathbf{b}$ . Thus, e.g.,  $\mathbf{abbabbaab}$  is contained in the language of  $\alpha$ , whereas  $\mathbf{bbbbaa}$  is not. It is a direct consequence of these definitions that a pattern language is either a singleton or infinite. Furthermore, it is worth noting that two basic types of pattern languages are considered in the literature, depending on whether the variables must stand for nonempty words (referred to as non erasing or NE-pattern languages) or whether they may represent the empty word (so-called extended, erasing or simply E-pattern languages).

While the definition of pattern languages is simple, many of their properties are known to be related to complex phenomena in combinatorics on words, such as pattern avoidability (see Jiang et al. [7]) and ambiguity of morphisms (see Reidenbach [12]). Hence, the knowledge on pattern languages is still patchy, despite recent progress mainly regarding decision problems (see, e. g., Freydenberger, Reidenbach [5], Fernau, Schmid [3], Fernau et al. [4] and Reidenbach, Schmid [13]) and the relation to the Chomsky hierarchy (see Jain et al. [6] and Reidenbach, Schmid [14]).

Establishing the closure properties of a class of formal languages is one of the most classical and fundamental research tasks in formal language theory and any respective progress normally leads to insights and techniques that yield a better understanding of the class. In the case of pattern languages, it is known since Angluin's initial work that they are not closed under most of the usual operations, including union, intersection and complement. However, these non-closure properties can be shown by using very basic example patterns and exploiting peculiarities of the definition of pattern languages. For example, if a pattern does not contain a variable, then its language is a singleton; hence the union of any two distinct singleton pattern languages contains two elements, and therefore it cannot be a pattern language. Furthermore, the intersection of two pattern languages given by patterns that start with different terminal symbols is empty and the empty set, although a trivial language, is not a pattern language as well. Since, apart from a strong result by Shinohara [15] on the union of NE-pattern languages, hardly anything is known beyond such immediate facts, we can observe that in the case of pattern languages the existing closure properties fail to contribute to our understanding of their intrinsic properties.

It is the main purpose of this paper to investigate the closure properties of pattern languages more thoroughly. To this end, in Section 3, we consider the closure properties of two important subclasses of pattern languages, namely the classes of terminal-free NE- and E-pattern languages, i. e., pattern languages that are generated by patterns that do not contain any terminal symbols. This choice is motivated by the fact that terminal-free patterns have been a recent focus of interest in the research on pattern languages and, furthermore, most existing examples for non-closure of pattern languages (including the two examples for union and intersection given in the previous paragraph) do not translate to the terminal-free case. In Section 3.1, we completely characterise when the union of two terminal-free pattern languages is again a terminal-free pattern language and, in Section 3.2, we prove their non-closure under intersection, for which the situation is much more complicated compared to the operation of union.

We consider general pattern languages in Section 4, and we provide complex examples demonstrating that it is probably a very hard task to obtain full characterisations of those pairs of pattern languages whose unions or intersections are again a pattern language. In Section 4.3, we also study the question whether an E-pattern language can be expressed by the union of *nonerasing* pattern languages and, likewise, whether an NE-pattern language can be expressed by the union of *erasing* pattern languages. This question is slightly at odds with the

classical investigation of closure properties, since we apply a language operation to members of one class and ask whether the resulting language is a member of another class. However, in the case of pattern languages, this makes sense, since every NE-pattern language is a finite union of E-pattern languages and every E-pattern language is a finite union of NE-pattern languages (see Jiang et al. [7]), a phenomenon that has been widely utilised in the context of inductive inference of pattern languages (see, e.g., Wright [17], Shinohara, Arimura [16]).

Due to space constraints, all proofs have been omitted from this paper.

## 2 Definitions and Preliminary Results

The symbols  $\cup$ ,  $\cap$  and  $\setminus$  denote the set operations of *union*, *intersection* and *set difference*, respectively. For sets  $U$  and  $B$  with  $B \subseteq U$ ,  $\overline{B} := U \setminus B$  is the *complement* of  $B$ .

Let  $\mathbb{N} := \{1, 2, 3, \dots\}$  and let  $\mathbb{N}_0 := \mathbb{N} \cup \{0\}$ . For an arbitrary alphabet  $A$ , a *word* (over  $A$ ) is a finite sequence of symbols from  $A$ , and  $\varepsilon$  stands for the *empty word*. The notation  $A^+$  denotes the set of all nonempty words over  $A$ , and  $A^* := A^+ \cup \{\varepsilon\}$ . For the *concatenation* of two words  $w_1, w_2$  we write  $w_1 \cdot w_2$  or simply  $w_1w_2$ , and  $w^n$  stands for the  $n$ -fold concatenation of the word  $w$ . We say that a word  $v \in A^*$  is a *factor* of a word  $w \in A^*$  if there are  $u_1, u_2 \in A^*$  such that  $w = u_1 \cdot v \cdot u_2$ . If  $u_1$  (or  $u_2$ ) is the empty word, then  $v$  is a *prefix* (or a *suffix*, respectively) of  $w$ . The notation  $|K|$  stands for the size of a set  $K$  or the length of a word  $K$ . A word  $w$  is *primitive* if, for any  $u$  such that  $w = u^k$ ,  $k = 1$ . The *primitive root* of a word  $w$  is the primitive word  $u$  such that  $w = u^k$ ,  $k \in \mathbb{N}$ .

For any alphabets  $A, B$ , a *morphism* is a function  $h : A^* \rightarrow B^*$  that satisfies  $h(vw) = h(v)h(w)$  for all  $v, w \in A^*$ ;  $h$  is said to be *nonerasing* if, for every  $a \in A$ ,  $h(a) \neq \varepsilon$ . A morphism  $h$  is *ambiguous* (with respect to a word  $w$ ) if there exists a morphism  $g$  satisfying  $g(w) = h(w)$  and, for a letter  $a$  in  $w$ ,  $g(a) \neq h(a)$ . If such a morphism  $g$  does not exist, then  $h$  is called *unambiguous* (with respect to  $w$ ). A morphism  $\sigma : A^* \rightarrow B^*$  is *periodic* if for some (primitive) word  $w \in B^*$ ,  $\sigma(x) \in \{w\}^*$  for every  $x \in A$ . The word  $w$  will be referred to as the *primitive root* of  $\sigma$ . If  $|\sigma(x)| = 1$  for every  $x \in A$ , then  $\sigma$  is *1-uniform*.

Let  $\Sigma$  be a finite alphabet of so-called *terminal symbols* and  $X$  a countably infinite set of *variables* with  $\Sigma \cap X = \emptyset$ . We normally assume  $X := \{x_1, x_2, x_3, \dots\}$ . A *pattern* is a nonempty word over  $\Sigma \cup X$ , a *terminal-free pattern* is a nonempty word over  $X$ ; if a word contains symbols from  $\Sigma$  only, then we occasionally call it a *terminal word*. For any pattern  $\alpha$ , we refer to the set of variables in  $\alpha$  as  $\text{var}(\alpha)$ . If the variables in a pattern  $\alpha$  are labelled in the natural way, then it is said to be in *canonical form*, i. e.,  $\alpha$  is in canonical form if, for some  $n \in \mathbb{N}$ ,  $\text{var}(\alpha) = \{x_1, x_2, \dots, x_n\}$  and, for any  $x_i, x_j \in \text{var}(\alpha)$  with  $i < j$ , there is a prefix  $\beta$  of  $\alpha$  such that  $x_i \in \text{var}(\beta)$  and  $x_j \notin \text{var}(\beta)$ . A pattern  $\alpha$  is a *one-variable pattern* if  $|\text{var}(\alpha)| = 1$ . A morphism  $h : (\Sigma \cup X)^* \rightarrow (\Sigma \cup X)^*$  is *terminal-preserving* if  $h(a) = a$  for every  $a \in \Sigma$ . The *residual* of a pattern  $\alpha$  is the word  $h_\varepsilon(\alpha)$ , where  $h_\varepsilon : (\Sigma \cup X)^* \rightarrow (\Sigma \cup X)^*$  is a terminal preserving morphism with  $h_\varepsilon(x) := \varepsilon$  for every  $x \in \text{var}(\alpha)$ . A terminal-preserving morphism  $h : (\Sigma \cup X)^* \rightarrow \Sigma^*$  is called a *substitution*.

**Definition 1.** Let  $\Sigma$  be an alphabet, and let  $\alpha \in (\Sigma \cup X)^*$  be a pattern. The E-pattern language of  $\alpha$  is defined by  $L_{E,\Sigma}(\alpha) := \{h(\alpha) \mid h : (\Sigma \cup X)^* \rightarrow \Sigma^* \text{ is a substitution}\}$ . The NE-pattern language of  $\alpha$  is defined by  $L_{NE,\Sigma}(\alpha) := \{h(\alpha) \mid h : (\Sigma \cup X)^* \rightarrow \Sigma^* \text{ is a nonerasing substitution}\}$ .

Note that we call a pattern language terminal-free if there exists a terminal-free pattern that generates it.

Some parts of our reasoning in the subsequent sections is based on word equations, which are defined as follows. For a set of unknowns  $Y$ , a terminal alphabet  $\Sigma$ , and two words  $\alpha, \beta \in (Y \cup \Sigma)^+$ , the expression  $\alpha = \beta$  is called a *word equation*. The solutions are terminal-preserving morphisms  $\sigma : (Y \cup \Sigma)^* \rightarrow \Sigma^*$  such that  $\sigma(\alpha) = \sigma(\beta)$ . The words  $\sigma(\alpha)$  ( $= \sigma(\beta)$ ) will be referred to as *solution-words*. It is often convenient to interpret variables from patterns as unknowns, and so word equations will often be formulated from two patterns.

This concludes the basic definitions of this paper. We now begin our investigation of the closure properties of the class of pattern languages. As a starting point, we refer to the corresponding result in the initial paper on pattern languages:

**Theorem 1 (Angluin [1]).** *NE-pattern languages are not closed under union, intersection, complement, Kleene plus, homomorphism and inverse homomorphism. NE-pattern languages are closed under concatenation and reversal.*

### 3 Terminal-Free Patterns

As briefly explained in Section 1, the proof of Theorem 1 heavily relies on the fact that patterns can contain terminal symbols. In the present section, we therefore wish to study whether the situation changes if we consider the classes of terminal-free E-pattern languages and terminal-free NE-pattern languages.

#### 3.1 Union

Simple examples show that neither the terminal-free NE-pattern languages nor the terminal-free E-pattern languages are closed under union:

**Proposition 1.** *Let  $\Sigma$  be an alphabet with  $\{a, b\} \subseteq \Sigma$ . For every  $Z, Z' \in \{E, NE\}$ , there does not exist a pattern  $\gamma$ , such that  $L_{Z,\Sigma}(\gamma) = L_{Z',\Sigma}(x_1x_1) \cup L_{Z',\Sigma}(x_1x_1x_1)$ .*

It is worth noting that the above statement also provides a first minor insight into the topic of expressing E-pattern languages as unions of NE-pattern languages and vice versa. We shall study this subject in Section 4.3 for patterns with terminal symbols in much more detail. In the present section, we merely want to point out that the union of two terminal-free E-pattern languages is indeed never a terminal-free NE-pattern language, and the union of two terminal-free NE-pattern languages cannot be a terminal-free E-pattern language:

**Proposition 2.** *Let  $\Sigma$  be an arbitrary alphabet, and let  $\alpha$  and  $\beta$  be terminal-free patterns. Then there does not exist a terminal-free pattern  $\gamma$  with  $L_{E,\Sigma}(\alpha) \cup L_{E,\Sigma}(\beta) = L_{NE,\Sigma}(\gamma)$  or  $L_{NE,\Sigma}(\alpha) \cup L_{NE,\Sigma}(\beta) = L_{E,\Sigma}(\gamma)$ .*

In the remainder of this section we wish to prove a similarly strong result for the actual closure of the class of terminal-free E- or NE-pattern languages. Hence, we wish to characterise those pairs of terminal-free (NE-/E-)pattern languages where the union again is a terminal-free (NE-/E-)pattern language. Our results shall demonstrate that the union of two terminal-free E-pattern languages can only be a terminal-free E-pattern language if there is an inclusion relation between the two languages, and that the same holds for the NE-pattern languages.

Our reasoning on the E case is based on a result on the inclusion problem for E-pattern languages. In [8], Jiang et al. provide a construction for a morphism  $\tau_k$  such that, for two patterns  $\alpha$  and  $\beta$ , the word  $\tau_{|\beta|}(\alpha)$  is contained in  $L_{E,\Sigma}(\beta)$  if and only if there exists a morphism  $\varphi$  from  $\beta$  to  $\alpha$ . This in turn implies that the erasing languages of two terminal free patterns satisfy a subset relation if and only if there exists a morphism from one pattern to the other. It is not difficult to see that this construction can be further used to satisfy, for patterns  $\alpha, \beta, \gamma$ , and  $k = \max(|\alpha|, |\beta|)$ , that  $\tau_k(\gamma) \in L_{E,\Sigma}(\alpha) \cup L_{E,\Sigma}(\beta)$  if and only if  $\gamma$  is a morphic image of  $\alpha$  or  $\beta$ . Thus if the relation  $L_{E,\Sigma}(\alpha) \cup L_{E,\Sigma}(\beta) = L_{E,\Sigma}(\gamma)$  is satisfied, then  $L_{E,\Sigma}(\gamma)$  is a subset of (and therefore also equal to)  $L_{E,\Sigma}(\alpha)$  or  $L_{E,\Sigma}(\beta)$  and we have the following situation.

**Lemma 1.** *Let  $\Sigma$  be an alphabet,  $|\Sigma| \geq 2$ , and let  $\alpha$  and  $\beta$  be terminal-free patterns. There exists a terminal-free pattern  $\gamma$  with  $L_{E,\Sigma}(\alpha) \cup L_{E,\Sigma}(\beta) = L_{E,\Sigma}(\gamma)$  if and only if  $L_{E,\Sigma}(\alpha) \subseteq L_{E,\Sigma}(\beta)$  or  $L_{E,\Sigma}(\beta) \subseteq L_{E,\Sigma}(\alpha)$ .*

It can be observed from simple examples that, in the nonerasing case, inclusion cannot be characterised by the existence of a morphism between the generating patterns. Thus, no equivalent argument can be derived for the nonerasing case. However, a corresponding result can be obtained by looking at the shortest words in the nonerasing languages of  $\alpha, \beta$  and  $\gamma$ . To this end, we define, for a pattern  $\alpha$ , the set  $M_\alpha$  to be  $\{\sigma(\alpha) \mid \sigma : \text{var}(\alpha)^* \rightarrow \Sigma^* \text{ is 1-uniform}\}$ .

The set  $M_\alpha$  has been used to positive effect in existing literature (see, e.g., [9]). It is particularly useful when considering nonerasing pattern languages because it encodes exactly the original pattern  $\alpha$  (up to a renaming of variables). Moreover, it has a number of convenient properties when considering the union of two NE-pattern languages. One such example is that if  $\alpha$  is strictly shorter than  $\beta$ , then the set of shortest words in  $L_{NE,\Sigma}(\alpha) \cup L_{NE,\Sigma}(\beta)$  will be exactly  $M_\alpha$ . Thus, if the union is itself the nonerasing language of some pattern  $\gamma$ , we have that  $\gamma = \alpha$  up to a renaming of variables. A similar result can be obtained for the case that  $|\alpha| = |\beta|$  by considering  $|M_\alpha \cup M_\beta|$ .

**Lemma 2.** *Let  $\Sigma$  be an alphabet,  $|\Sigma| \geq 2$ , and let  $\alpha, \beta$  be terminal free patterns in canonical form with  $|\alpha| = |\beta|$ . Suppose that  $\gamma$  is a terminal free pattern (again in canonical form) with  $M_\alpha \cup M_\beta = M_\gamma$ . Then  $\gamma \in \{\alpha, \beta\}$ .*

Consequently, we can verify the same statement for nonerasing languages as we have for erasing languages.

**Lemma 3.** *Let  $\Sigma$  be an alphabet,  $|\Sigma| \geq 2$ , and let  $\alpha$  and  $\beta$  be terminal-free patterns. There exists a terminal-free pattern  $\gamma$  with  $L_{NE,\Sigma}(\alpha) \cup L_{NE,\Sigma}(\beta) = L_{NE,\Sigma}(\gamma)$  if and only if  $L_{NE,\Sigma}(\alpha) \subseteq L_{NE,\Sigma}(\beta)$  or  $L_{NE,\Sigma}(\beta) \subseteq L_{NE,\Sigma}(\alpha)$ .*

Note that Lemma 3 extends an equivalent result by Shinohara [15] that holds for alphabets with at least 3 letters.

Thus, in general, the languages of two terminal-free patterns only union together to produce a third in the trivial case.

**Theorem 2.** *Let  $Z, Z' \in \{E, NE\}$ . Let  $\Sigma$  be an alphabet,  $|\Sigma| \geq 2$ , and let  $\alpha, \beta, \gamma$  be terminal-free patterns. Then  $L_{Z,\Sigma}(\alpha) \cup L_{Z,\Sigma}(\beta) = L_{Z',\Sigma}(\gamma)$  if and only if  $L_{Z,\Sigma}(\alpha) = L_{Z',\Sigma}(\gamma)$  and  $L_{Z,\Sigma}(\beta) \subseteq L_{Z,\Sigma}(\alpha)$  or  $L_{Z,\Sigma}(\beta) = L_{Z',\Sigma}(\gamma)$  and  $L_{Z,\Sigma}(\alpha) \subseteq L_{Z,\Sigma}(\beta)$ .*

It is worth noting that, for terminal-free patterns, the inclusion problem – and therefore the question of closure under union – is decidable in the E case (see Jiang et al. [8], as explained above), but still open in the NE case.

### 3.2 Intersection

In the present section, we wish to investigate if the terminal-free NE- or E-pattern languages are closed under intersection. For the NE case, simple counterexamples such as  $\alpha := xyx$  and  $\beta := xxy$  can be used to prove the following observation:

**Proposition 3.** *The terminal-free NE-pattern languages are not closed under intersection.*

We can obtain an equivalent result for the terminal-free E-pattern languages, but our reasoning is significantly more complex and requires the analysis of certain word equations. Moreover, we are able to provide a characterisation for a restricted class of pairs of patterns, and show that, for this class, the situation is non-trivial (i.e., there exist both positive and negative examples). We proceed by considering the link between word equations and intersections of pattern-languages.

If, for a word equation  $\alpha = \beta$ , the words  $\alpha$  and  $\beta$  are over disjoint alphabets, then the set of solutions  $\sigma : (\text{var}(\alpha) \cup \text{var}(\beta))^* \rightarrow \Sigma^*$  corresponds exactly to the set of pairs of morphisms  $\tau_1 : \text{var}(\alpha)^* \rightarrow \Sigma^*$ ,  $\tau_2 : \text{var}(\beta)^* \rightarrow \Sigma^*$  such that  $\tau_1(\alpha) = \tau_2(\beta)$ . Thus, it also exactly describes the intersection  $L_{E,\Sigma}(\alpha) \cap L_{E,\Sigma}(\beta)$ . Furthermore, such an intersection is invariant under renamings of  $\alpha$  and of  $\beta$ , so any intersection of E-pattern languages can be described in this way. The next proposition gives a characterisation of when the intersection of two terminal-free E-pattern languages is again a terminal-free E-pattern language in the restricted case that the corresponding word equation permits only periodic solutions. Note that, for  $\alpha$  and  $\beta$  over disjoint alphabets, such solutions always exist.

**Proposition 4.** *Let  $\Sigma$  be an arbitrary alphabet. Let  $\alpha, \beta$  be terminal-free patterns over disjoint alphabets and suppose that the word equation  $\alpha = \beta$  permits only periodic solutions. Let  $w$  be the shortest non-empty solution-word. Let*

$$\mu := \text{lcm}(\text{gcd}\{|\alpha|_{x_i} \mid x_i \in \text{var}(\alpha)\}, \text{gcd}\{|\beta|_{y_j} \mid y_j \in \text{var}(\beta)\}).$$

*Then  $L_{E,\Sigma}(\alpha) \cap L_{E,\Sigma}(\beta)$  is a terminal-free E-pattern language if and only if  $\mu = |w|$ .*

Despite Proposition 4, it is still a non-trivial task to find two terminal-free E-pattern languages whose intersection is not a terminal-free E-pattern language. In particular, it remains to find appropriate patterns  $\alpha$  and  $\beta$  such that the word equation  $\alpha = \beta$  has only periodic solutions. The following proposition provides such an example, and hence we have the analogous result to Proposition 3.

**Proposition 5.** *Let  $\Sigma$  be an arbitrary alphabet, and let  $\alpha := x_1x_2x_1^2x_2x_1^3x_2^2$  and  $\beta := x_3x_4^2x_3^2x_4^6x_3^3$ . Then  $L_{E,\Sigma}(\alpha) \cap L_{E,\Sigma}(\beta)$  cannot be expressed as a terminal-free E-pattern language.*

It is even possible to give a much stronger statement, showing the extent to which the ‘pattern-language mechanism’ is incapable of handling this seemingly uncomplicated set of solutions.

**Corollary 1.** *For any alphabet  $\Sigma$ ,  $L_{E,\Sigma}(x_1x_2x_1^2x_2x_1^3x_2^2) \cap L_{E,\Sigma}(x_3x_4^2x_3^2x_4^6x_3^3)$  cannot be expressed as a finite union of terminal-free E-pattern languages.*

It is worth noting that the approach above can be used to show that for  $\alpha' := x_1x_2x_1^2x_2^2x_1^3x_2^3$  and  $\beta' := x_3x_4^2x_3^2x_4^7x_3^3$ , one has that  $L_{E,\Sigma}(\alpha') \cap L_{E,\Sigma}(\beta') = L_{E,\Sigma}(x_1^6)$ . This demonstrates that the intersection of two E-pattern languages can in some cases be expressed as an E-pattern language, and therefore that the problem of whether the intersection of two E-pattern languages form an E-pattern language is nontrivial. However it is worth pointing out that a characterisation of this situation is probably very difficult to acquire due to the challenging nature of finding solution-sets of word equations.

### 3.3 Other Closure Properties

In this Section, we show that regarding the closure under the operations of complementation, morphisms, inverse morphisms, Kleene plus and Kleene star, terminal-free pattern languages behave similarly to the full class of pattern languages.

**Proposition 6.** *For every terminal-free pattern  $\alpha$ ,  $\overline{L_{E,\Sigma}(\alpha)}$  is not a terminal-free E-pattern language and  $\overline{L_{NE,\Sigma}(\alpha)}$  is not a terminal-free NE-pattern language.*

Proposition 6 does not only prove the non-closure of terminal-free E- and NE-pattern languages under complementation, but also characterises in a trivial way the terminal-free pattern languages whose complement is also a terminal-free pattern language.

**Proposition 7.** *Let  $\Sigma$  be a terminal alphabet with  $|\Sigma| \geq 2$ . The terminal-free NE- and E-pattern languages, with respect to  $\Sigma$ , are not closed under morphisms, inverse morphisms, Kleene plus and Kleene star.*

## 4 General Patterns

As explained in Section 1 and formally stated in Theorem 1, the closure properties of the full classes of NE-pattern languages and of E-pattern languages are understood. In the present section, we therefore wish to expand the more specific insights into the terminal-free pattern languages gained in Section 3 to the full classes. More precisely, with respect to the operations of complementation, intersection and union, we investigate those patterns that exhibit the property that their complement, intersection or union is again a pattern language and we try to characterise these patterns. Our strongest results are with respect to the operation of union.

### 4.1 Complement

With respect to the full class of E- and NE-pattern language, an analogue of Proposition 6 exists:

**Proposition 8 (Bayer [2]).** *Let  $\Sigma$  be a terminal alphabet with  $|\Sigma| \geq 2$ . For every pattern  $\alpha$ ,  $L_{E,\Sigma}(\alpha)$  is not an E-pattern language and  $L_{NE,\Sigma}(\alpha)$  is not an NE-pattern language.*

In the same way as Proposition 6 does for terminal-free patterns, this proposition yields a trivial characterisation of pattern languages with a complement that again is a pattern language.

### 4.2 Intersection

It is straightforward to construct patterns  $\alpha$  and  $\beta$  such that  $L_{E,\Sigma}(\alpha) \cap L_{E,\Sigma}(\beta)$  is not an E-pattern language or  $L_{NE,\Sigma}(\alpha) \cap L_{NE,\Sigma}(\beta)$  is not an NE-pattern language. Furthermore, any two terminal-free patterns  $\alpha$  and  $\beta$  are an example for the situation that  $L_{E,\Sigma}(\alpha) \cap L_{E,\Sigma}(\beta)$  is not an NE-pattern language and, as long as there are at least two symbols in  $\Sigma$ , also for the situation that  $L_{NE,\Sigma}(\alpha) \cap L_{NE,\Sigma}(\beta)$  is not an E-pattern language. Moreover, there are non-trivial examples of patterns  $\alpha$ ,  $\beta$  and  $\gamma$ , such that  $L_{NE,\Sigma}(\alpha) \cap L_{NE,\Sigma}(\beta) = L_{E,\Sigma}(\gamma)$ :

- $L_{NE,\Sigma}(\mathbf{ax}) \cap L_{NE,\Sigma}(xx) = L_{E,\Sigma}(\mathbf{axax})$ .
- $L_{NE,\Sigma}(x\mathbf{ay}) \cap L_{NE,\Sigma}(xxx) = L_{E,\Sigma}(x\mathbf{ayxayxay})$ .
- $L_{NE,\Sigma}(\mathbf{axa}) \cap L_{NE,\Sigma}(xx) = L_{E,\Sigma}(\mathbf{axaaxa})$ .
- $L_{NE,\Sigma}(\mathbf{axax}) \cap L_{NE,\Sigma}(x\mathbf{bxb}) = L_{E,\Sigma}(\mathbf{axbaxb})$ .
- $L_{NE,\Sigma}(\mathbf{axy}) \cap L_{NE,\Sigma}(xxx) = L_{E,\Sigma}(\mathbf{axaxax})$ .

However, it is not known whether or not there are patterns  $\alpha$  and  $\beta$ , such that  $L_{E,\Sigma}(\alpha) \cap L_{E,\Sigma}(\beta)$  is an NE-pattern language. Moreover, we do not have any characterisations for the situation that the intersection of two pattern languages is again a pattern language.



### 4.3 Union

Examples of patterns  $\alpha$  and  $\beta$  such that  $L_{Z,\Sigma}(\alpha) \cup L_{Z,\Sigma}(\beta)$  is not a  $Z'$ -pattern language, for all  $Z, Z' \in \{E, NE\}$ , are provided by Proposition 1.

Theorem 2 is our strongest result in Section 3, as it shows that the union of terminal-free pattern languages can only be a terminal-free pattern language if one of the languages is contained in the other. At first glance it seems a reasonable hypothesis that a similar result might hold for the full class of pattern languages, but in the present section we show that this is not true.

For all but the union of pairs of E-pattern languages and the question of whether they can form an E-pattern language, suitable examples are not too hard to find:

**Proposition 9.** *Let  $\Sigma$  be a terminal alphabet.*

- $L_{E,\{a,b\}}(\mathbf{aax}) \cup L_{E,\{a,b\}}(\mathbf{abx}) = L_{NE,\{a,b\}}(\mathbf{ax})$ .
- $L_{NE,\Sigma}(\mathbf{abc}) \cup L_{NE,\Sigma}(\mathbf{axbxcx}) = L_{E,\Sigma}(\mathbf{axbxcx})$ .
- $L_{NE,\{a,b\}}(\mathbf{ax_1}) \cup L_{NE,\{a,b\}}(\mathbf{bx_1}) = L_{NE,\{a,b\}}(x_1x_2)$ .

Regarding the question of whether  $L_{E,\Sigma}(\alpha) \cup L_{E,\Sigma}(\beta) = L_{E,\Sigma}(\gamma)$  for patterns  $\alpha, \beta, \gamma$  implies that there is an inclusion relation between  $L_{E,\Sigma}(\alpha)$  and  $L_{E,\Sigma}(\beta)$ , the following three propositions provide increasingly complex counterexamples for alphabet sizes 2, 3, and 4.

**Proposition 10.** *Let  $\Sigma = \{a, b\}$ ,  $\alpha := x_1\mathbf{ax_2bx_2ax_3}$ ,  $\beta := x_1\mathbf{ax_2bbx_2ax_3}$  and  $\gamma := x_1\mathbf{ax_2bx_3ax_4}$ . Then  $L_{E,\Sigma}(\alpha) \not\subseteq L_{E,\Sigma}(\beta)$ ,  $L_{E,\Sigma}(\beta) \not\subseteq L_{E,\Sigma}(\alpha)$  and  $L_{E,\Sigma}(\alpha) \cup L_{E,\Sigma}(\beta) = L_{E,\Sigma}(\gamma)$ .*

**Proposition 11.** *Let  $\Sigma := \{a, b, c\}$ ,*

$$\begin{aligned} \alpha &:= x_1\mathbf{ax_2x_3^6x_4^3x_5^6x_6bx_7ax_2x_8^{12}x_9^6x_{10}^{12}x_6bx_{10}}, \\ \beta &:= x_1\mathbf{ax_2x_3^6x_4^2x_5^5x_6^6x_7bx_8ax_2x_9^{12}x_{10}^4x_{11}^{10}x_{12}^{10}x_7bx_{11}} \text{ and} \\ \gamma &:= x_1\mathbf{ax_2x_3^6x_4^2x_5^3x_6^6x_7bx_8ax_2x_9^{12}x_{10}^4x_{11}^6x_{12}^{10}x_7bx_{11}}. \end{aligned}$$

*Then*

$$L_{E,\Sigma}(\alpha) \not\subseteq L_{E,\Sigma}(\beta), \quad L_{E,\Sigma}(\beta) \not\subseteq L_{E,\Sigma}(\alpha) \text{ and } L_{E,\Sigma}(\alpha) \cup L_{E,\Sigma}(\beta) = L_{E,\Sigma}(\gamma).$$

**Proposition 12.** *Let  $\Sigma := \{a, b, c, d\}$ ,*

$$\begin{aligned} \alpha &:= x_1\mathbf{ax_2x_3^2x_4^2x_5^2x_6bx_7ax_2x_8^2x_9^2x_{10}^2x_6bx_{10}cx_{11}x_{12}^2x_{13}^2x_{14}^2x_{15}^2x_{16}^2d} \\ &\quad x_{17}cx_{11}x_{18}^2x_{19}^2x_{20}^2x_{21}^2x_{22}^2x_{23}^2x_{24}^2x_{25}^2x_{26}^2x_{27}^2x_{28}^2x_{29}^2x_{30}^2x_{31}^2x_{32}^2x_{33}^2x_{34}^2x_{35}^2x_{36}^2, \\ \beta &:= x_1\mathbf{ax_2x_3^2x_4^2x_5^2x_6^2x_7bx_8ax_2x_9^2x_{10}^2x_{11}^2x_{12}^2x_{13}^2x_{14}^2x_{15}^2x_{16}^2d} \\ &\quad x_{17}cx_{12}x_{18}^2x_{19}^2x_{20}^2x_{21}^2x_{22}^2x_{23}^2x_{24}^2x_{25}^2x_{26}^2x_{27}^2x_{28}^2x_{29}^2x_{30}^2x_{31}^2x_{32}^2x_{33}^2x_{34}^2x_{35}^2 \text{ and} \\ \gamma &:= x_1\mathbf{ax_2x_3^2x_4^2x_5^2x_6^2x_7bx_8ax_2x_9^2x_{10}^2x_{11}^2x_{12}^2x_{13}^2x_{14}^2x_{15}^2x_{16}^2x_{17}^2d} \\ &\quad x_{18}cx_{12}x_{19}^2x_{20}^2x_{21}^2x_{22}^2x_{23}^2x_{24}^2x_{25}^2x_{26}^2x_{27}^2x_{28}^2x_{29}^2x_{30}^2x_{31}^2x_{32}^2x_{33}^2x_{34}^2x_{35}^2. \end{aligned}$$

*Then*

$$L_{E,\Sigma}(\alpha) \not\subseteq L_{E,\Sigma}(\beta), \quad L_{E,\Sigma}(\beta) \not\subseteq L_{E,\Sigma}(\alpha) \text{ and } L_{E,\Sigma}(\alpha) \cup L_{E,\Sigma}(\beta) = L_{E,\Sigma}(\gamma).$$

We are not able to give equivalent examples for larger alphabets, and we expect the question of their existence to be a complex and important problem. This is because the above examples depend on the ambiguity of terminal-preserving morphisms, which is a phenomenon that underpins many properties of pattern languages. Similar constructions to those in Propositions 10, 11, and 12 have been used to disprove longstanding conjectures on inductive inference (see Reidenbach [10,12]) of and the equivalence problem (see Reidenbach [11]) for E-pattern languages over alphabets of up to 4 letters and, similarly, it has so far not been possible to expand those techniques to arbitrary alphabets. Our examples, thus, suggest a close link between the problem in the current section and the two most important open problems for E-pattern languages over alphabets with at least 5 letters, and we expect that substantial progress on any one of them will require combinatorial insights that will allow the others to be solved as well.

For all  $Z, Z' \in \{E, NE\}$ , we have seen example patterns  $\alpha$  and  $\beta$  such that  $L_{Z,\Sigma}(\alpha) \cup L_{Z,\Sigma}(\beta)$  is a  $Z'$ -pattern language. We shall now try to generalise these examples in order to obtain characterisations of such pairs of patterns.

For the case  $Z = Z' = E$ , we are only able to state a necessary condition for  $L_{E,\Sigma}(\alpha) \cup L_{E,\Sigma}(\beta) = L_{E,\Sigma}(\gamma)$  that, unfortunately, is not very strong:

**Theorem 3.** *Let  $\Sigma$  be an alphabet, and let  $\alpha, \beta$  and  $\gamma$  be patterns with  $L_{E,\Sigma}(\alpha) \cup L_{E,\Sigma}(\beta) = L_{E,\Sigma}(\gamma)$ . Furthermore, let  $w_\alpha, w_\beta$  and  $w_\gamma$  be the residuals of  $\alpha, \beta$  and  $\gamma$ , respectively. Then  $w_\gamma = w_\alpha$  and  $w_\gamma$  is a subsequence of  $w_\beta$  or  $w_\gamma = w_\beta$  and  $w_\gamma$  is a subsequence of  $w_\alpha$ .*

In view of the fact that the examples of Propositions 10, 11 and 12 are rather complicated, we expect that a full characterisation for the case  $Z = Z' = E$  is difficult to obtain.

For the case  $Z = Z' = NE$ , we can present a strong necessary condition that, similarly to Lemma 3, strengthens a result by Shinohara [15]:

**Theorem 4.** *Let  $\Sigma$  be an alphabet with  $\{a, b\} \subseteq \Sigma$  and let  $\alpha, \beta$  and  $\gamma$  be patterns. If  $L_{NE,\Sigma}(\alpha) \cup L_{NE,\Sigma}(\beta) = L_{NE,\Sigma}(\gamma)$ , then one of the following three statements is true:*

- $L_{NE,\Sigma}(\alpha) \subseteq L_{NE,\Sigma}(\beta)$  and  $\beta = \gamma$ .
- $L_{NE,\Sigma}(\beta) \subseteq L_{NE,\Sigma}(\alpha)$  and  $\alpha = \gamma$ .
- $|\Sigma| = 2$  and

$$\begin{aligned} \alpha &= \delta_0 \mathbf{a} \delta_1 \mathbf{a} \delta_2 \dots \delta_{m-1} \mathbf{a} \delta_m, \\ \beta &= \delta_0 \mathbf{b} \delta_1 \mathbf{b} \delta_2 \dots \delta_{m-1} \mathbf{b} \delta_m, \\ \gamma &= \delta_0 x \delta_1 x \delta_2 \dots \delta_{m-1} x \delta_m, \end{aligned}$$

where  $m \geq 1, \delta_i \in (X \cup \Sigma)^*, 0 \leq i \leq m$ .

It remains to consider the cases  $Z = NE, Z' = E$  and  $Z = E, Z' = NE$ , for which we have full characterisations. Before we prove these characterisations, we recall that Jiang et al. show in [7] that, for every pattern  $\alpha$ , we can construct finite

sets of patterns  $\Gamma$  and  $\Delta$  such that  $L_{E,\Sigma}(\alpha) = \bigcup_{\beta \in \Gamma} L_{NE,\Sigma}(\beta)$  and  $L_{NE,\Sigma}(\alpha) = \bigcup_{\beta \in \Delta} L_{E,\Sigma}(\beta)$ . More precisely,  $\Gamma$  is the set of all patterns that can be obtained from  $\alpha$  by erasing some (possibly none) of the variables and  $\Delta$  contains all pattern that can be obtained from  $\alpha$  by substituting each  $x \in \text{var}(\alpha)$  by  $bx$ , for some  $b \in \Sigma$ . We note that the examples  $L_{NE,\Sigma}(\mathbf{abc}) \cup L_{NE,\Sigma}(\mathbf{axbxcx}) = L_{E,\Sigma}(\mathbf{axbxcx})$  and  $L_{E,\{a,b\}}(\mathbf{aax}) \cup L_{E,\{a,b\}}(\mathbf{abx}) = L_{NE,\{a,b\}}(\mathbf{ax})$  of Proposition 9 are applications of exactly this construction.

The characterisation for the case  $Z = NE, Z' = E$  follows from the fact that we can prove that if we restrict ourselves to unions of only two pattern languages, then  $L_{E,\Sigma}(\alpha) = \bigcup_{\beta \in \Gamma} L_{NE,\Sigma}(\beta)$  is the only possible way to describe an E-pattern language by NE-pattern languages.

**Theorem 5.** *Let  $\Sigma$  be an alphabet with  $|\Sigma| \geq 2$  and let  $\alpha, \beta$  and  $\gamma$  be patterns. Then  $L_{NE,\Sigma}(\alpha) \cup L_{NE,\Sigma}(\beta) = L_{E,\Sigma}(\gamma)$  if and only if  $\alpha \in \Sigma^+$  and  $\beta = \gamma = u_1 x^{j_1} u_2 x^{j_2} \dots x^{j_m} u_{m+1}$ ,  $j_i \in \mathbb{N}_0$ ,  $1 \leq i \leq m$ , such that  $u_1 u_2 \dots u_{m+1} = \alpha$ .*

With respect to the case  $Z = E, Z' = NE$ , we can even present a characterisation for the situation  $L_{NE,\Sigma}(\alpha) = \bigcup_{i=1}^k L_{E,\Sigma}(\beta_i)$  with  $k \leq |\Sigma|$ . It shall be explained later on that this characterisation is a generalisation of the construction given by Jiang et al.

**Theorem 6.** *Let  $\ell \geq 2$  and let  $\Sigma$  be an alphabet with  $\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_\ell\} \subseteq \Sigma$ . Furthermore, let  $\alpha_1, \alpha_2, \dots, \alpha_\ell$  and  $\gamma$  be patterns with  $L_{E,\Sigma}(\alpha_i) \neq L_{E,\Sigma}(\alpha_j)$ ,  $1 \leq i < j \leq \ell$ . Then  $\bigcup_{i=1}^\ell L_{E,\Sigma}(\alpha_i) = L_{NE,\Sigma}(\gamma)$  if and only if, for some permutation  $\pi$  of  $(1, 2, \dots, \ell)$ ,*

- $\Sigma = \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_\ell\}$ ,
- $\gamma = u_1 x u_2 x u_3 \dots u_k x u_{k+1}$ ,  $k \geq 1$ ,  $u_i \in \Sigma^*$ ,  $1 \leq i \leq k + 1$ , and,
- for every  $i$ ,  $1 \leq i \leq \ell$ ,

$$\alpha_i = u_1 \alpha'_i \mathbf{a}_{\pi(i)} \alpha''_i u_2 \alpha'_i \mathbf{a}_{\pi(i)} \alpha''_i u_3 \dots u_k \alpha'_i \mathbf{a}_{\pi(i)} \alpha''_i u_{k+1},$$

where  $\alpha'_i, \alpha''_i \in X^*$ ,

- for every  $i$ ,  $1 \leq i \leq \ell$ , there exists a  $y_i \in \text{var}(\alpha_i)$  with  $|\alpha_i|_{y_i} = k$  and
  - $|\alpha'_i|_{y_i} = 1$  for all  $i$ ,  $1 \leq i \leq \ell$ , or
  - $|\alpha''_i|_{y_i} = 1$  for all  $i$ ,  $1 \leq i \leq \ell$ .

If we apply the construction of Jiang et al. to a one-variable pattern  $\gamma$ , then we obtain patterns  $\alpha_i$ ,  $1 \leq i \leq |\Sigma|$ , that satisfy the conditions of the patterns in the statement of Theorem 6. More precisely, this corresponds to the special case where  $\alpha'_i \alpha''_i = y_i$ ,  $1 \leq i \leq |\Sigma|$ . Moreover, it can be easily verified that if  $\gamma$  and patterns  $\alpha_i$ ,  $1 \leq i \leq |\Sigma|$ , satisfy the conditions of the statement of Theorem 6, then, depending on whether  $|\alpha'_i|_{y_i} = 1$  for all  $i$ ,  $1 \leq i \leq |\Sigma|$ , or  $|\alpha''_i|_{y_i} = 1$  for all  $i$ ,  $1 \leq i \leq |\Sigma|$ , we can obtain patterns  $\beta_i$  from the patterns  $\alpha_i$  by replacing  $\alpha'_i \mathbf{a}_i \alpha''_i$  by  $y_i \mathbf{a}_i$  or by  $\mathbf{a}_i y_i$ , respectively, and  $\bigcup_{i=1}^{|\Sigma|} L_{E,\Sigma}(\beta_i) = L_{NE,\Sigma}(\gamma)$  still holds. Furthermore, the patterns  $\beta_i$  are exactly the patterns that are obtained if we apply the construction of Jiang et al.

**Acknowledgments.** The authors wish to thank the anonymous referees for their helpful suggestions, which have yielded a stronger version of Proposition 7.

## References

1. Angluin, D.: Finding patterns common to a set of strings. *Journal of Computer and System Sciences* 21, 46–62 (1980)
2. Bayer, H.: Allgemeine Eigenschaften von Patternsprachen. Projektarbeit, Fachbereich Informatik, Universität Kaiserslautern (2007) (in German)
3. Fernau, H., Schmid, M.L.: Pattern matching with variables: A multivariate complexity analysis. In: Fischer, J., Sanders, P. (eds.) *CPM 2013. LNCS*, vol. 7922, pp. 83–94. Springer, Heidelberg (2013)
4. Fernau, H., Schmid, M.L., Villanger, Y.: On the parameterised complexity of string morphism problems. In: *Proc. 33rd IARCS Annual Conference on Foundations of Software Technology and Theoretical Computer Science, FSTTCS 2013. Leibniz International Proceedings in Informatics (LIPIcs)*, vol. 24, pp. 55–66 (2013)
5. Freydenberger, D., Reidenbach, D.: Bad news on decision problems for patterns. *Information and Computation* 208, 83–96 (2010)
6. Jain, S., Ong, Y., Stephan, F.: Regular patterns, regular languages and context-free languages. *Information Processing Letters* 110, 1114–1119 (2010)
7. Jiang, T., Kinber, E., Salomaa, A., Salomaa, K., Yu, S.: Pattern languages with and without erasing. *International Journal of Computer Mathematics* 50, 147–163 (1994)
8. Jiang, T., Salomaa, A., Salomaa, K., Yu, S.: Decision problems for patterns. *Journal of Computer and System Sciences* 50, 53–63 (1995)
9. Lange, S., Wiehagen, R.: Polynomial-time inference of arbitrary pattern languages. *New Generation Computing* 8, 361–370 (1991)
10. Reidenbach, D.: A non-learnable class of E-pattern languages. *Theoretical Computer Science* 350, 91–102 (2006)
11. Reidenbach, D.: An examination of Ohlebusch and Ukkonen’s conjecture on the equivalence problem for E-pattern languages. *Journal of Automata, Languages and Combinatorics* 12, 407–426 (2007)
12. Reidenbach, D.: Discontinuities in pattern inference. *Theoretical Computer Science* 397, 166–193 (2008)
13. Reidenbach, D., Schmid, M.L.: Patterns with bounded treewidth. *Information and Computation* (to appear)
14. Reidenbach, D., Schmid, M.L.: Regular and context-free pattern languages over small alphabets. *Theoretical Computer Science* 518, 80–95 (2014)
15. Shinohara, T.: Inferring unions of two pattern languages. *Bulletin of Informatics and Cybernetics* 20, 83–88 (1983)
16. Shinohara, T., Arimura, H.: Inductive inference of unbounded unions of pattern languages from positive data. *Theoretical Computer Science* 241, 191–209 (2000)
17. Wright, K.: Identification of unions of languages drawn from an identifiable class. In: *Proc. 2nd Annual Workshop on Computational Learning Theory, COLT 1989*, pp. 328–333 (1989)