



# Excessively duplicating patterns represent non-regular languages



Carles Creus<sup>1,2</sup>, Guillem Godoy<sup>\*,1</sup>, Lander Ramos<sup>1</sup>

Universitat Politècnica de Catalunya, Omega building, Jordi Girona 1-3, Barcelona 08034, Spain

## ARTICLE INFO

### Article history:

Received 5 March 2013

Received in revised form 10 October 2013

Accepted 11 November 2013

Available online 13 November 2013

Communicated by A. Muscholl

### Keywords:

Theory of computation

Pattern

Regular tree language

Tree automaton

Tree homomorphism

## ABSTRACT

A constrained term pattern  $s : \varphi$  represents the language of all instances of the term  $s$  satisfying the constraint  $\varphi$ . For each variable in  $s$ , this constraint specifies the language of its allowed substitutions. Regularity of languages represented by sets of patterns has been studied for a long time. This problem is known to be co-NP-complete when the constraints allow each variable to be replaced by any term over a fixed signature, and EXPTIME-complete when the constraints restrict each variable to a regular set. In both cases, duplication of variables in the terms of the patterns is a necessary condition for non-regularity. This is because duplications force the recognizer to test equality between subterms. Hence, for the specific classes of constraints mentioned above, if all patterns are linear, then the represented language is necessarily regular. In this paper we focus on the opposite case, that is when there are patterns with “excessively duplicating” variables. We prove that when each pattern of a non-empty set has a duplicated variable constrained to an infinite language, then the language represented by the set is necessarily non-regular. We prove this result independently of the kind of constraints used, just assuming that they are mappings from variables to arbitrary languages. Our result provides an efficient procedure for detecting, in some cases, non-regularity of images of regular languages under tree homomorphisms.

© 2013 Elsevier B.V. All rights reserved.

## 1. Introduction

A constrained term pattern (or just a pattern) is a pair  $s : \varphi$ , where  $s$  is a term and  $\varphi$  maps each variable occurring in  $s$  to a language. The pattern  $s : \varphi$  represents the language of all terms obtained from  $s$  by replacing each variable  $x$  occurring in  $s$  by a term in  $\varphi(x)$ . Thus,  $\varphi$  constrains the instances of  $s$  by restricting the possible substitutions of each variable. A set of patterns  $S$  represents the union of the languages represented by each pattern in  $S$ .

Patterns are a widely used formalism in computer science to represent languages. As with many other representation formalisms, one is frequently interested in solving questions like whether a given term belongs to a language, or whether a language is a subset of another language. Another important question is to determine whether the language represented by a set of patterns is (tree) regular, i.e., recognizable by a tree automaton [1]. This question is relevant because regular tree languages have a more tractable representation by means of tree automata, and have better properties like closure by several set operations, and existence of efficient procedures for some of such operations and for several decision problems.

The expressive power and decidable properties of patterns depend on the mechanism used to describe the constraints. For example, for unrestricted patterns, i.e., when we allow to replace each variable by any term over a fixed signature, regularity has been proved co-NP-complete [2,3].

\* Corresponding author. Tel.: +34 93 4137815; fax: +34 93 4137833.

E-mail addresses: ccreuslopez@gmail.com (C. Creus), ggodoy@lsi.upc.edu (G. Godoy), landertxu@gmail.com (L. Ramos).

<sup>1</sup> The three authors were supported by Spanish Ministry of Education and Science by the FORMALISM project (TIN2007-66523).

<sup>2</sup> Supported by an FPU grant from the Spanish Ministry of Education.

The problem is EXPTIME-complete when the constraints restrict each variable to a regular set [4].

Patterns with variables occurring at least twice in the term and constrained to an infinite language are often the cause of non-regularity. For example, a pattern of the form  $f(x, x) : \varphi$ , where  $\varphi$  maps  $x$  to an infinite language, represents a non-regular set, since tree automata only have a finite number of states and cannot test equality between subterms.

In contrast to patterns with duplicated variables, linear patterns, i.e., the ones where each variable occurs at most once in the term, increase the chance of regularity. For example, for unrestricted patterns and for patterns with regular constraints, if all patterns are linear, then the represented language is necessarily regular. However, when considering more expressive constraints, linearity of the patterns does not necessarily imply regularity of the represented language.

In this paper we focus on proving non-regularity when the patterns of a given set have “excessively duplicating” variables. More precisely, we prove that when each pattern of a set has a duplicated variable constrained to an infinite language, then the language represented by the set is necessarily non-regular. This property holds independently of the kind of constraints used, i.e., we just assume that they are mappings from variables to arbitrary languages.

Finally, we apply our result to show an efficient (linear time) algorithm that detects, in some cases, non-regularity of images of regular languages under tree homomorphisms. Note that the general case has recently been proved EXPTIME-complete [5].

The paper is structured as follows. In Section 2 we present basic preliminary definitions used in the rest of the paper. In Section 3 we define patterns, the particular case of excessively duplicating patterns, and prove that sets of excessively duplicating patterns always represent non-regular languages. In Section 4 we apply our result on patterns to detect, in some cases, the non-regularity of the image of a regular language under a tree homomorphism. We conclude in Section 5.

## 2. Preliminaries

### 2.1. Terms

We use the standard notation from the term rewriting literature [6]. The powerset of a set  $S$  is denoted by  $2^S$ , and the cardinal of  $S$  is denoted by  $\#S$ . A signature  $\Sigma$  is a (finite) set of function symbols with arity, which is partitioned as  $\bigcup_i \Sigma^{(i)}$  such that  $f \in \Sigma^{(m)}$  if the arity of  $f$  is  $m$ . We sometimes denote  $\Sigma$  explicitly as  $\{f_1 : m_1, \dots, f_n : m_n\}$ , where  $f_1, \dots, f_n$  are the function symbols and  $m_1, \dots, m_n$  are the corresponding arities. Symbols in  $\Sigma^{(0)}$ , called *constants*, are denoted by  $a, b$ , with possible subscripts. The elements of a set  $\mathcal{X}$  of variable symbols are denoted by  $x, y, z$  with possible subscripts. The set  $\mathcal{T}(\Sigma, \mathcal{X})$  of *terms* over  $\Sigma$  and  $\mathcal{X}$ , is the smallest set containing  $\mathcal{X}$  such that  $f(t_1, \dots, t_m)$  is in  $\mathcal{T}(\Sigma, \mathcal{X})$  whenever  $f \in \Sigma^{(m)}$  and  $t_1, \dots, t_m \in \mathcal{T}(\Sigma, \mathcal{X})$ . The set of variables occurring in a term  $t$  is denoted  $\text{vars}(t)$ . A term  $t$  is called *linear* if each variable occurs at most once in  $t$ . A term  $t$  is called *ground*

if it contains no variables. The set of all ground terms over  $\Sigma$  is denoted  $\mathcal{T}(\Sigma)$ . A *language* over  $\Sigma$  is a set of ground terms.

A *position* is a sequence of natural numbers. The symbol  $\lambda$  denotes the empty sequence, and  $p.p'$  denotes the concatenation of the positions  $p$  and  $p'$ . The set of positions of a term  $t$ , denoted  $\text{Pos}(t)$ , is defined recursively as  $\text{Pos}(f(t_1, \dots, t_m)) = \{\lambda\} \cup \{i.p \mid i \in \{1, \dots, m\} \wedge p \in \text{Pos}(t_i)\}$ . The *length* of a position is denoted  $|p|$ . Note that  $|\lambda| = 0$  and  $|i.p| = 1 + |p|$ . A position  $p_1$  is a *prefix* of a position  $p$ , denoted  $p_1 \leq p$ , if there is a position  $p_2$  such that  $p_1.p_2 = p$ . Also,  $p_1$  is a *proper prefix* of  $p$ , denoted  $p_1 < p$ , if  $p_1 \leq p$  and  $p_1 \neq p$ . Two positions  $p, p'$  are *parallel*, denoted  $p \parallel p'$ , if  $p \not\leq p'$  and  $p' \not\leq p$ .

The *subterm* of a term  $t$  at a position  $p \in \text{Pos}(t)$ , denoted  $t|_p$ , is defined recursively as  $t|_\lambda = t$  and  $f(t_1, \dots, t_m)|_{i.p} = t_i|_p$ . The *replacement* of the subterm at position  $p$  in a term  $t$  by a term  $s$ , denoted  $t[s]_p$ , is defined recursively as  $t[s]_\lambda = s$  and  $f(t_1, \dots, t_m)[s]_{i.p} = f(t_1, \dots, t_i[s]_p, \dots, t_m)$ .

The *root* of a term  $t = f(t_1, \dots, t_m)$  is  $\text{root}(t) = f$ . A term  $t$  over  $\Sigma$  can be seen as a function from its set of positions into  $\Sigma$ . For this reason, we shall denote by  $t(p)$  the symbol labeling  $t$  at position  $p$ , i.e.,  $t(p) = \text{root}(t|_p)$ .

The *height* of a term  $t$ , denoted  $\text{height}(t)$ , is defined recursively as  $\text{height}(t) = 0$  if  $t$  is a variable or a constant, and as  $\text{height}(f(t_1, \dots, t_m)) = 1 + \max(\text{height}(t_1), \dots, \text{height}(t_m))$  otherwise.

A *substitution*  $\sigma$  is a mapping from variables to terms. It can also be used as a function from terms to terms recursively defined by  $\sigma(f(t_1, \dots, t_m)) = f(\sigma(t_1), \dots, \sigma(t_m))$  for terms different from variables. For example, if  $\sigma$  is  $\{x \mapsto f(b, y), y \mapsto a\}$ , then  $\sigma(g(x, y))$  is  $g(f(b, y), a)$ .

### 2.2. Automata

A *tree automaton* (TA) is a tuple  $A = \langle Q, \Sigma, F, \Delta \rangle$ , where  $Q$  is a finite set of states,  $\Sigma$  is a signature,  $F \subseteq Q$  is the subset of final states (also called accepting states), and  $\Delta$  is a set of rules of the form  $f(q_1, \dots, q_m) \rightarrow q$ , where  $q_1, \dots, q_m, q \in Q$  and  $f \in \Sigma^{(m)}$ . The *size* of  $A$  is defined by  $\#A = \#Q$ .

A *run*  $r$  of a TA  $A$  on a term  $t \in \mathcal{T}(\Sigma)$  is a function  $r : \text{Pos}(t) \rightarrow Q$  satisfying that, for each position  $p \in \text{Pos}(t)$ , if  $t|_p$  is of the form  $f(t_1, \dots, t_m)$ , then there exists a rule of the form  $f(q_1, \dots, q_m) \rightarrow q$  in  $\Delta$  such that  $r(p.1) = q_1, \dots, r(p.m) = q_m$  and  $r(p) = q$ . A run  $r$  is *accepting* if  $r(\lambda)$  is accepting. A term  $t$  is *accepted* or *recognized* by  $A$  if there exists an accepting run of  $A$  on  $t$ . The language recognized by  $A$ , denoted  $\mathcal{L}(A)$ , is the set of terms accepted by  $A$ . By  $\mathcal{L}(A, q)$  we denote the set of terms for which there exists a run  $r$  of  $A$  such that  $r(\lambda) = q$ .

We say that a language  $L$  is *regular* if there exists a TA  $A$  such that  $\mathcal{L}(A) = L$ .

Given a TA  $A$ , a term  $t$  and a run  $r$  of  $A$  on  $t$ , we define  $r|_p$  as the run of  $A$  on  $t|_p$  described by  $r|_p(p') = r(p.p')$ . In addition, given a run  $r'$  of  $A$  on a term  $t'$  satisfying  $r'(\lambda) = r(p)$ , we define  $r[r']_p$  as the run of  $A$  on  $t[t']_p$  described by  $r[r']_p(p') = r(p')$  for positions  $p'$  holding  $p \not\leq p'$ , and otherwise  $r[r']_p(p.p') = r'(p')$ .

Download English Version:

<https://daneshyari.com/en/article/427413>

Download Persian Version:

<https://daneshyari.com/article/427413>

[Daneshyari.com](https://daneshyari.com)