



Prime normal form and equivalence of simple grammars

Cédric Bastien^a, Jurek Czyzowicz^a, Wojciech Fraczak^{a, b, *}, Wojciech Rytter^c

^a*Dépt d'informatique, Université du Québec en Outaouais, Gatineau PQ, Canada*

^b*IDT Canada Inc., Ottawa, Ont., Canada*

^c*Institute of Informatics, Warsaw University, Warsaw, Poland*

Abstract

A prefix-free language is prime if it cannot be decomposed into a concatenation of two prefix-free languages. We show that we can check in polynomial time if a language generated by a simple context-free grammar is prime. Our algorithm computes a canonical representation of a simple language, converting its arbitrary simple grammar into prime normal form (PNF); a simple grammar is in PNF if all its nonterminals define primes. We also improve the complexity of testing the equivalence of simple grammars. The best previously known algorithm for this problem worked in $O(n^{13})$ time. We improve it to $O(n^7 \log^2 n)$ and $O(n^5 \text{polylog } v)$ time, where n is the total size of the grammars involved, and v is the length of a shortest string derivable from a nonterminal, maximized over all nonterminals.

© 2006 Elsevier B.V. All rights reserved.

Keywords: Deterministic context-free language; Language equivalence; Simple grammar; Prime normal form

1. Introduction

An important question in language theory is, given a class of languages, how to find a canonical representation of any language of this class. Such a representation often permits solving various decidability problems related to a given class of languages, such as the equality of languages, nonemptiness, etc. Most often the canonical representation of the language is given by a special form of its grammar, called a normal form. In this paper, we give an algorithm converting a simple grammar into its equivalent, unique representation in a form of so-called prime normal form (PNF). The canonical form of simple grammar was studied by Courcelle, c.f. [2]. The crucial question that our algorithm is confronted with, is whether a simple language is prime, i.e., not decomposable into a concatenation of two nontrivial prefix-free languages.

In general, the canonical representation of any type of language may be substantially larger than its original grammar. This is also the case for simple languages. Hence verifying the equivalence of simple languages by means of canonical representations may be inefficient. The equivalence problem for simple context-free grammars is a classical question in formal language theory. It is a nontrivial problem, since the inclusion problem for simple languages is undecidable. Korenjak and Hopcroft, see [11,8], proved that the equivalence problem is decidable and they gave the first, doubly

* Corresponding author. Institute of Informatics, Warsaw University, Warsaw, Poland.

E-mail address: fraczak@gmail.com (W. Fraczak).

exponential time algorithm solving it. Their result was improved by Caucal to $O(n^3 v(G))$ time, see [1]. The parameter n is the size of the simple grammar and $v(G)$ is the length of a shortest string derived from a nonterminal, maximized over all nonterminals. Caucal's algorithm is exponential since $v(G)$ can be exponential with respect to n . Hirshfeld, Jerrum, and Moller gave the first polynomial $O(n^{13})$ time algorithm for this problem in [9]. We call it the HJM algorithm.

In the second part of the paper, we design an algorithm based on a version of Caucal's algorithm, that has a better complexity than HJM. More precisely, our algorithm works in time $O(n^7 \log^2 n)$. On the other hand, a variation of our algorithm works in time $O(n^5 \text{polylog}(v(G)))$, thus beating the complexity of Caucal's algorithm, e.g., for $v(G) \in \Omega(n^3)$. Similarly as the HJM algorithm, we apply the techniques used in the algorithmic theory of compressed strings, based on Lempel–Ziv string encoding. The idea of such an encoding is that, instead of representing a string explicitly, we design a context-free grammar generating the string as a one-word language. As the combinatorial complexity of such grammar can be significantly smaller than the length of the word, it may be considered as a succinct representation of the word. Such encodings were recently considered by researchers, mainly in the context of efficient pattern matching. There is one problem in this field which is of particular interest to us—the *compressed first mismatch problem* (*First-MP*). Given two strings encoded by a grammar, *First-MP* looks for the position of the first symbol at which the strings differ. Polynomial time algorithms for computing *First-MP* were given independently in [9,13], in very disjoint settings. More powerful algorithms were proposed in [10], where a more complicated problem of fully compressed string-matching was solved. For the purpose of this paper, we will use the technique from [12], which we adopted to obtain a faster algorithm.

Simple languages are applied by IDT Canada to perform packet classification at wire speed. Classes of packets are described with the aid of simple languages, and their recognition is made by a so-called Concatenation State Machine, an efficient version of a stateless pushdown automaton. As shown in [5], there is a one-to-one correspondence between Concatenation State Machines and simple grammars. In order to store large sets of classification policies in memory, it is necessary to reuse their common parts. A natural way to do this consists of decomposing simple languages into primes, each of which is stored in memory only once. When a new classification policy is added to memory, we verify if its prime factors are already stored in the data base. The algorithms described in this paper are used to decompose classification policies into primes and to identify primes for reuse.

2. Simple languages

A context-free grammar $G = (\Sigma, N, P)$ is composed of a finite set Σ of *terminals*, a finite set N of *nonterminals* disjoint from Σ , and a finite set $P \subset N \times (N \cup \Sigma)^*$ of *production rules*. For every $\beta, \gamma \in (N \cup \Sigma)^*$, if $(A, \alpha) \in P$, then $\beta A \gamma \rightarrow \beta \alpha \gamma$. A *derivation* $\beta \xrightarrow{*} \gamma$ is a finite sequence $(\alpha_0, \alpha_1, \dots, \alpha_n)$ such that $\beta = \alpha_0$, $\gamma = \alpha_n$, and $\alpha_{i-1} \rightarrow \alpha_i$ for $i \in [1, n]$.

For every sequence of nonterminals $\alpha \in N^*$ of a grammar $G = (\Sigma, N, P)$, we denote by $L_G(\alpha)$ the set of terminal strings derivable from α , i.e., $L_G(\alpha) \stackrel{\text{def}}{=} \{w \in \Sigma^* \mid \alpha \xrightarrow{*} w\}$. Often, if G is known from the context, we will write $L(\alpha)$ instead of $L_G(\alpha)$.

A grammar $G = (\Sigma, N, P)$ is in *Greibach normal form* if for every production rule $(A \rightarrow \alpha) \in P$, we have $\alpha \in \Sigma N^*$. A grammar $G = (\Sigma, N, P)$ is a *simple context-free grammar* (simple grammar) if G is a Greibach normal form grammar and such that whenever $A \rightarrow a\alpha_1$ and $A \rightarrow a\alpha_2$, for a same $a \in \Sigma$, then $\alpha_1 = \alpha_2$.

A language $L \subseteq \Sigma^*$ is a *simple language* (also called *s-language*) if $L = \{\varepsilon\}$ (where ε denotes the empty word) or if there exists a simple grammar $G = (\Sigma, N, P)$ such that $L_G(A) = L$, for some $A \in N$. The definition implies that every nonterminal of a simple grammar defines a simple language. Since simple languages are prefix codes and are closed by concatenation, the family of simple languages under concatenation constitutes a free monoid with $\{\varepsilon\}$ as unit. Thus, every nontrivial simple language L (i.e., $L \neq \{\varepsilon\}$ and $L \neq \emptyset$) admits a unique decomposition into *prime* (i.e., undecomposable, nontrivial) simple languages, $L = P_1 P_2 \dots P_n$.

3. PNF for simple grammars

In this section we give an algorithm converting any simple grammar to its canonical representation called *PNF*. A simple grammar is in PNF if each of its nonterminals represents a prime. We will use the following algebraic notation

for left and right division in the free monoid of prefix codes. If $L = L_1 L_2$ for some prefix codes L, L_1, L_2 , then by $L_1^{-1} L$ we denote L_2 and by $L L_2^{-1}$ we denote L_1 . We call L_1 a left divider and L_2 a right divider of L .

Let L be a prefix code and $L = P_1 P_2 \dots P_n$ be its decomposition into primes. Prime P_n will be called *final prime* of L , and it will be denoted by $f(L)$. In particular, if L is a prime, then $f(L) = L$.

Lemma 1. *Let $G = (\Sigma, N, P)$ be a simple grammar. For every $X \in N$, there exists $Y \in N$, such that $f(L(X)) = L(Y)$.*

Proof. Let $w \in L(X) f(L(X))^{-1}$, and $X \xrightarrow{*} w\alpha$ be the leftmost derivation in G , with $\alpha \in N^+$. Since $L(\alpha) = f(L(X))$ and $L(\alpha)$ is a prime, α consists of a single nonterminal, i.e., $\alpha \in N$. \square

Let $w_0 \alpha_0 \rightarrow \dots \rightarrow w_i \alpha_i \rightarrow \dots \rightarrow w_n \alpha_n$ be the leftmost derivation $X \xrightarrow{*} w$, with $w_0 = \varepsilon$, $\alpha_0 = X$, $w_n = w$, $\alpha_n = \varepsilon$, $w_i \in \Sigma^*$, and $\alpha_i \in N^*$, for $i \in [0, n]$. We are interested in the subsequence $\pi(X, w) = Y_0, Y_1, \dots, Y_j$ of $\alpha_0, \alpha_1, \dots, \alpha_n$, which consists of those elements of $\alpha_0, \dots, \alpha_n$ that are single nonterminals. E.g., for the leftmost derivation of $abcdef \in L(X)$:

$$\underline{X} \rightarrow aY\underline{Y} \rightarrow ab\underline{Y} \rightarrow abc\underline{Y} \rightarrow abcdYZ \rightarrow abcde\underline{Z} \rightarrow abcdef$$

we have $\pi(X, abcdef) = X, Y, Y, Z$.

Definition 2. Let $G = (\Sigma, N, P)$ be a simple grammar. We define relation \mathcal{D} over $N \cup \{\varepsilon\}$ as follows. $(X, Y) \in \mathcal{D}$ if and only if:

- there exists a rule $(X \rightarrow a\alpha Y)$ in P for some $a \in \Sigma$ and $\alpha \in N^*$ or
- $Y = \varepsilon$ and there exists a rule $(X \rightarrow a)$ in P for some $a \in \Sigma$.

Relation \mathcal{D} can be seen as a digraph $(N \cup \{\varepsilon\}, \mathcal{D}, \varepsilon)$ with sink ε . In a digraph with a sink, vertex v is called a *d-articulation point* of vertex u if and only if v is present on every path from u to the sink. It was shown in [4] that the order of first (or last) occurrences of the d-articulation points of a vertex v is the same in all paths from v to the sink. Thus, it is natural to represent the set of all d-articulation points for a given vertex v as an ordered list of vertices, (u_0, u_1, \dots, u_n) , where $u_0 = v$ and u_n is the sink.

In [4], it was shown that a prefix code L is prime if and only if the initial state v_1 of the minimal deterministic automaton for L does not have any d-articulation point except sink and v_1 itself. Moreover, the list of d-articulation points (v_1, v_2, \dots, v_n) corresponds to the prime decomposition of L , the factors being the languages defined by automata having v_i as the initial state and v_{i+1} as the final state (with all outgoing transitions of the final state removed), for $i \in [1, n]$, respectively.

Lemma 3. *For every path π from X to ε in \mathcal{D} there exists a word $w \in L(X)$, such that $\pi = \pi(X, w)$. Conversely, for every $w \in L(X)$, $\pi(X, w)$ defines a path from X to ε in \mathcal{D} .*

We say that a grammar $G = (\Sigma, N, P)$ is *reduced* if there is no two different nonterminals defining the same language, i.e., for all $X, Y \in N$, if $L(X) = L(Y)$ then $X = Y$. By Lemma 1, the set of nonterminals $F(X) \stackrel{\text{def}}{=} \{Y \in N \mid L(Y) = f(L(X))\}$ is nonempty. If the underlying grammar is reduced then $F(X)$ consists of a single nonterminal which, by convenient abuse of notation, will be denoted by $f(X)$.

Theorem 4. *Let $G = (\Sigma, N, P)$ be a reduced simple grammar. For every $X \in N$, $L(X)$ is prime if and only if X does not have d-articulation points in \mathcal{D} except sink and X itself. Moreover, if $Y \in N$ is a d-articulation point of X then $L(Y)$ is a right divider of $L(X)$.*

Proof. By Lemma 1, since G is reduced, every derivation starting in X is of form $X \xrightarrow{*} w' f(X) \xrightarrow{*} w$. Thus, for every $w \in L(X)$, $\pi(X, w)$ contains $f(X)$, i.e., $f(X)$ is a d-articulation point of X in \mathcal{D} .

Let Y be a d-articulation point of X in \mathcal{D} . By Lemma 3, every derivation starting in X passes by Y , thus Y is a d-articulation point for X in the (infinite) deterministic automaton for X , which implies that $L(Y)$ is a right divider of $L(X)$, cf. [6]. \square

Input: Simple grammar $G = (\Sigma, N, P)$ and $S \in N^+$.

Output: Simple grammar G' in PNF and S' such that $L_G(S) = L_{G'}(S')$.

(1) Reduce G .

Find redundant nonterminals by checking if $L(X) = L(Y)$, for all $X, Y \in N$.
The redundant nonterminals are substituted in P and S , and removed from N .

(2) For every $X \in N$, find $f(X) \in N$.

Construct \mathcal{D} for G and find the second-last d-articulation point for X .
If for every $X \in N$, $X = f(X)$, then **return** (G, S) .

(3) Construct a new grammar $G' = (\Sigma, N, P')$ and new S' :

Define morphism $h : N \mapsto N^*$ as: $h(X) \stackrel{\text{def}}{=} \begin{cases} X & \text{if } X = f(X) \\ Xf(X) & \text{otherwise.} \end{cases}$

Set S' to $h(S)$, and P' as follows, for $a \in \Sigma$, $X, Y \in N$, $\alpha \in N^*$:

- (a) If $(X \rightarrow a\alpha) \in P$ and $X = f(X)$, then $(X \rightarrow ah(\alpha))$ is in P' .
- (b) If $(X \rightarrow a\alpha f(X)) \in P$ and $X \neq f(X)$, then $(X \rightarrow ah(\alpha))$ is in P' .
- (c) If $(X \rightarrow a\alpha Y) \in P$, $X \neq f(X)$, $Y \neq f(X)$, then $(X \rightarrow ah(\alpha)Y)$ is in P' .

(4) Set G to G' , S to S' and go to 1.

Fig. 1. Algorithm $PNF(G, S)$.

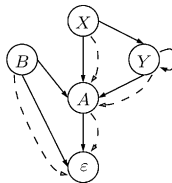
Theorem 5. Given a reduced simple grammar $G = (\Sigma, N, P)$, we can find $f(X)$ for all $X \in N$ in linear time.

Proof. By Theorem 4, the nonterminal $f(X)$ is exactly the second last d-articulation point for X in \mathcal{D} . Calculating $f(X)$ for all $X \in N$ can be done in linear time, by using an algorithm for finding dominators in flow graphs, cf. [7]. \square

The algorithm for transforming a simple grammar $G = (\Sigma, N, P)$ into PNF, called $PNF(G, S)$, is presented in Fig. 1.

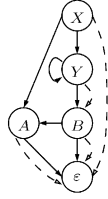
Example 6. We present an example of the execution of the algorithm. The input consists of a simple grammar $G = \{(X \rightarrow aAA), (X \rightarrow bYY), (Y \rightarrow aY), (Y \rightarrow bBA), (A \rightarrow a), (B \rightarrow aXA), (B \rightarrow b)\}$, and a simple language represented as a word $S = XA$ over nonterminals of G . We obtain the grammar in PNF while keeping track of the decomposition of S . For each iteration, we give the value of S , the grammar G , the digraph \mathcal{D} (solid edges), the d-articulation tree (broken line edges), and the values $f(x)$ for $x \in \{X, Y, A, B\}$ and $h(x)$ for $x \in \{X, Y, A, B, S\}$.

$S = XA$
 $X = aAA + bYY$
Iteration 1: $Y = aY + bBA$
 $A = a$
 $B = aXA + b$



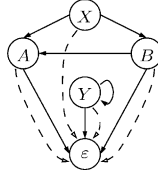
x	$f(x)$	$h(x)$
X	A	XA
Y	A	YA
A	A	A
B	B	B
S	$-$	XAA

$$\begin{aligned}
S &= XAA \\
X &= aA + bYAY \\
\text{Iteration 2: } Y &= aY + bB \\
A &= a \\
B &= aXAA + b
\end{aligned}$$



x	$f(x)$	$h(x)$
X	X	X
Y	B	YB
A	A	A
B	B	B
S	–	XAA

$$\begin{aligned}
S &= XAA \\
X &= aA + bYBAYB \\
\text{Iteration 3: } Y &= aY + b \\
A &= a \\
B &= aXAA + b
\end{aligned}$$



x	$f(x)$	$h(x)$
X	X	X
Y	Y	Y
A	A	A
B	B	B
S	–	XAA

Theorem 7. *The algorithm $PNF(G, S)$ correctly computes a PNF simple grammar G' and S' such that $L_G(S) = L_{G'}(S')$.*

Proof. Step 1 does not change the semantics of any nonterminal, so it reduces G to an equivalent simple grammar. Step 2 effectively finds final primes for all nonterminals. Step 3 transforms the grammar G into G' by right-factorizing every nonprime nonterminal X by $f(X)$: if X is prime then $L_G(X) = L_{G'}(X)$, otherwise $L_G(X) = L_{G'}(Xf(X))$. Every production $(X \rightarrow \alpha) \in P$ is rewritten accordingly into a corresponding production $(X \rightarrow \beta) \in P'$. Hence, for all $X \in N$, $L_G(X) = L_{G'}(h(X))$. Thus morphism h converts grammar G together with S to a grammar G' with $S' = h(S)$ such that $L_G(S) = L_{G'}(S')$. Every iteration of the program cuts the length of nonprime nonterminals, in terms of their prime decomposition, by one. Thus, the total number of iterations equals the maximum length of the prime decompositions of nonterminals of the initial grammar. Hence the algorithm terminates. By the exit condition from Step 2, each nonterminal is prime, hence G is in PNF. \square

Both steps 2 and 3, of the algorithm may be computed in linear time, hence the complexity of each iteration of the main loop is dominated by the grammar reduction from step 1.

The polynomial time algorithm from Section 5, repeated $O(n^2)$ times may be used to perform the grammar reduction. However, for grammar $G = \{(A_1 \rightarrow aA_2A_2), (A_2 \rightarrow aA_3A_3), \dots, (A_{n-1} \rightarrow aA_nA_n), (A_n \rightarrow a)\}$, language $L_G(A_1)$ has an exponential number of primes with respect to the size of G . Hence the number of iterations of the main loop of $PNF(G, S)$ may be exponential and so may be the size of the resulting PNF grammar. Since simple languages constitute a free monoid, the PNF is unique.

Corollary 8. *Every simple language L can be represented by a PNF simple grammar $G = (\Sigma, N, P)$ and a starting word $S \in N^*$, such that $L_G(S) = L$. Such a representation is unique. The problem of constructing the PNF representation of L given by a simple grammar is decidable. The PNF representation may be of exponential size with respect to the size of the original grammar.*

4. First mismatch-pair problem

Our approach to transform Caucal's algorithm for the equivalence problem of simple grammars, cf. [1], into a polynomial time one (with respect to the size of the input grammars) uses compressed representations of sequences of nonterminals, instead of using explicit representations.

We will use the terminology of *acyclic morphisms* because it is more convenient in presenting our algorithms. It is basically equivalent to the representation of a single word by a context-free grammar generating exactly one word, or to a “straight line program”.

A *morphism* over a monoid M is an application $H : M \mapsto M$ such that $H(1_M) = 1_M$ and $H(x \cdot y) = H(x) \cdot H(y)$, for all $x, y \in M$. A morphism $H : M \mapsto M$ is fully defined by providing the values for the generators of M . Thus, a morphism H over a finitely generated free monoid N^* is usually defined by providing $H : N \mapsto N^*$. A morphism $H : N \rightarrow N^+$ is said to be *acyclic* if we can order elements of N in such a way that for each $A \in N$, we have: $H(A) = A$ or $A > B$ for each symbol B occurring in the string $H(A)$. For an acyclic morphism H over N^* we denote $H^{|N|}$ by H^* , since $H^{|N|+1} = H^{|N|}$. If $H^*(\alpha) = w$ then we say that (H, α) is a compressed representation of w . The size of w can be exponential with respect to the size of its compressed representation.

Let $G = (\Sigma, N, P)$ be a simple grammar. We say that an acyclic morphism $H : N \mapsto N^+$ is *self-proving* in G if for each $A \in N$, we have

- if $A \rightarrow a\alpha$ then $H(A) \rightarrow a\beta$ and $H^*(\alpha) = H^*(\beta)$; and
- if $H(A) \rightarrow a\beta$ then $A \rightarrow a\alpha$ and $H^*(\alpha) = H^*(\beta)$.

The concept of self-proving relations was introduced by Courcelle, c.f. [3]. In order to show how to apply this idea, we need to introduce first the concept of *bisimulation*. A simple grammar $G = (\Sigma, N, P)$ can be seen as a compact description of a deterministic (infinite) automaton, where states are words over N and transitions are defined by production rules: $\alpha \xrightarrow{a} \beta$ if $\alpha = A\alpha'$, $(A \rightarrow a\gamma) \in P$, and $\beta = \gamma\alpha'$.

A relation $\mathcal{R} \subseteq N^* \times N^*$ is a *bisimulation* if and only if, whenever $(\alpha, \beta) \in \mathcal{R}$, we have

- if $\alpha \xrightarrow{a} \alpha'$ then $\beta \xrightarrow{a} \beta'$ for some β' such that $(\alpha', \beta') \in \mathcal{R}$; and
- if $\beta \xrightarrow{a} \beta'$ then $\alpha \xrightarrow{a} \alpha'$ for some α' such that $(\alpha', \beta') \in \mathcal{R}$.

We write $\alpha \sim \beta$ if $(\alpha, \beta) \in \mathcal{R}$ for some bisimulation \mathcal{R} . We write $(\alpha, \beta) \xrightarrow{a} (\alpha', \beta')$ to say that $\alpha \xrightarrow{a} \alpha'$ and $\beta \xrightarrow{a} \beta'$.

Lemma 9. Let $G = (\Sigma, N, P)$ be a simple grammar and $\alpha, \beta \in N^*$. $L_G(\alpha) = L_G(\beta)$ if and only if $\alpha \sim \beta$.

Proof. Both implications are easily provable by contradiction. \square

The following lemma reformulates the idea of Courcelle in the terms of acyclic morphisms.

Lemma 10. If H is an acyclic morphism self-proving in $G = (\Sigma, N, P)$, then $L_G(x) = L_G(H(x))$, for every $x \in N^*$.

Proof. Firstly, we prove that if H is self-proving and $(x, H^i(x)) \xrightarrow{a} (\alpha, \beta)$, for $x, \alpha, \beta \in N^+$ and $a \in \Sigma$, then $H^*(\alpha) = H^*(\beta)$, for any $i \geq 0$.

For $i = 0$, $H^i(x) = x$, and the implication is obviously true.

Let $A \in N, x, y, z \in N^*$, and $(Ax, H(Ax)) \xrightarrow{a} (yx, zH(x))$. From the definition of *self-proving* morphism, we have

$$H^*(yx) = H^*(y)H^*(x) = H^*(z)H^*(x) = H^*(z)H^*(H(x)) = H^*(zH(x)).$$

Let $(x, H^{k+n}(x)) \xrightarrow{a} (\alpha, \beta)$. By induction, for some $\gamma \in N^*$, $(x, H^k(x)) \xrightarrow{a} (\alpha, \gamma)$, $(H^k(x), H^n(H^k(x))) \xrightarrow{a} (\gamma, \beta)$, and $H^*(\alpha) = H^*(\gamma) = H^*(\beta)$.

Secondly, we prove that for every $x, y, x', y' \in N^*$, if $H^*(x) = H^*(y)$ and $(x, y) \xrightarrow{a} (x', y')$, then $H^*(x') = H^*(y')$. Since $H^*(x) = H^*(y)$, an existence of $(x, H^*(x)) \xrightarrow{a} (x', z)$ implies an existence of $(y, H^*(y)) \xrightarrow{a} (y', z)$. By the intermediate result proved above we have $H^*(x') = H^*(y')$.

Finally, we notice that relation $B \stackrel{\text{def}}{=} \{(x, y) \in N^* \times N^* \mid H^*(x) = H^*(y)\}$, contains all pairs $(x, H^*(x))$ and it is a bisimulation. Thus, by Lemma 9, $L_G(x) = L_G(H(x))$. \square

The crucial tool in the polynomial-time algorithms is the compressed *first mismatch-pair* problem, *First-MP*:

Input: An acyclic morphism $H : N \mapsto N^+$ and two strings $x, y \in N^*$.

Output:

- *First-MP*(x, y, H) = *nil*, if $H^*(x) = H^*(y)$;
- *First-MP*(x, y, H) = *failure*, if one of $H^*(x), H^*(y)$ is a proper prefix of the other;

- $\text{First-MP}(x, y, H) = (A, B) \in N \times N$, where (A, B) is the *first mismatch pair*, i.e., the first symbols occurring at the same position in $H^*(x)$ and in $H^*(y)$, respectively, which are different.

We say that a morphism H is binary if $|H(A)| \leq 2$ for each $A \in N$. The following fact has been essentially shown in [12].

Lemma 11. *Assume that a given acyclic morphism $H : N \mapsto N^+$ is binary and that the length of x and y is at most $O(|N|)$, then we can solve $\text{First-MP}(x, y, H)$ in time $O(k^2 \cdot h^2)$, where $k = |N|$ and $h \stackrel{\text{def}}{=} \min\{k \geq 0 \mid H^k = H^{k+1}\}$ is the depth of the morphism.*

5. The equivalence algorithm

Conceptually it is easier to deal with grammars in binary Greibach normal form (denoted GNF2). This means that each side of the production is of the form $(A \rightarrow a\alpha)$, where $a \in \Sigma$ and $\alpha \in \{\varepsilon\} \cup N \cup N^2$.

Lemma 12. *For each simple grammar G of total size n (the total number of symbols describing G) there is an equivalent simple grammar G' in GNF2 with only $O(n)$ nonterminal symbols. G' can be constructed from G in $O(n)$ time.*

The total size of a grammar in GNF2 is of a same order as the size of N . Hence by the size of a grammar $G = (\Sigma, N, P)$, we mean $n = |N|$.

All known algorithms for the equivalence problem for simple grammars are based on the possibility of computing the quotient of one prefix language by another, assuming that the quotient exists and the languages are given as two nonterminals of a simple grammar.

More precisely, let A and B be two nonterminals of a simple grammar $G = (\Sigma, N, P)$, such that $L(A) = L(B) \cdot L$, for some language $L \subseteq \Sigma^*$. The language L can be derived from A by a leftmost derivation following any word w from $L(B)$, i.e., $A \xrightarrow{*} w\gamma$, for $\gamma \in N^*$, and $L(\gamma) = L$.

Let $\|A\|$ denote the length of the shortest word derivable from A .

Lemma 13. *Let G be a simple grammar of size n . We can compute the lengths of shortest words derivable from all nonterminals of G in time $O(n \log n)$.*

Proof. Finding $\|A\|$ for all $A \in N$ corresponds to the single-source shortest paths problem in an *and/or* graph, which, using Dijkstra algorithm, can be solved in time $O(n \log n)$. \square

Lemma 14. *Let A and B be two nonterminals of a simple grammar $G = (\Sigma, N, P)$ such that $L(A) = L(B) \cdot L$, for some $L \subseteq \Sigma^*$. We can compute $\gamma \in N^*$ such that $L(\gamma) = L$ in time $O(n)$. It is guaranteed that $|\gamma| \leq n$.*

Proof. Consider the parse tree for the derivation of a shortest word w from A . The idea is to find a path down the tree which cuts off left of this path subtrees γ generating prefix of w of length $\|B\|$. Since w is a shortest word, no path of the parse tree contains two occurrences of the same nonterminal hence the depth of the tree is at most n . Therefore, $|\gamma| \leq n$ and computing takes $O(n)$ time. \square

The result of the algorithm for finding the quotient of A by B as described in the proof of Lemma 14 will be denoted by $\text{quot}(A, B)$. The algorithm will give a result for any pair of nonterminals A and B , as long as $\|A\| \geq \|B\|$. Notice that $L(A) = L(B \text{quot}(A, B))$ only if $L(B)$ is a left divider of $L(A)$.

Lemma 10 is the starting point for the design of the algorithm EQUIVALENCE. Assume that we fix a linear order $A_1 < A_2 < \dots < A_n$ of nonterminals, such that whenever $i < j$, we have $\|A_i\| \leq \|A_j\|$. The idea of the algorithm is to construct a self-proving morphism H or, in the process of its construction, to discover a failure which contradicts $L(A) = L(B)$. The main point of the algorithm is to keep pairs of long strings in compressed form. We keep only strings of linear length, their explicit representations are determined by the morphism H . Each time a new rule is generated by setting $H(A) = B\gamma$, where $\gamma = \text{quot}(A, B)$, we create pairs (α, β) such that $A \rightarrow a\alpha$ and $B\gamma \rightarrow a\beta$, for every letter a of the terminal alphabet. We keep the generated pairs in set Q . To each pair we apply operation *First-MP*, which “eliminates” the next nonterminal, or finds that we have a pair of identical strings, such pairs are removed from Q .

```

Input: Simple grammar  $G = (\Sigma, N, P)$  and nonterminals  $X, Y \in N$ ;
Output: TRUE if  $L_G(X) = L_G(Y)$ , FALSE otherwise.

Initialization:
   $Q := \{(X, Y)\}$ ;
  for each  $A \in N$  do  $H[A] := A$ ;
  while  $Q$  is not empty do
     $(\beta_1, \beta_2) :=$  an element of  $Q$ ;
    switch ( $First-MP(\beta_1, \beta_2, H)$ ) do
      case nil : remove  $(\beta_1, \beta_2)$  from  $Q$ ;
      case failure : return FALSE;
      case  $(A, B)$  :
         $\gamma := quot(A, B)$ ;
         $H[A] := B\gamma$ ; /* The nonterminal  $A$  is "eliminated" */
        for each  $a \in \Sigma$  do
          if  $(A, B\gamma) \xrightarrow{a} (\beta'_1, \beta'_2)$  then add  $(\beta'_1, \beta'_2)$  to  $Q$ ;
          if  $(A \xrightarrow{a} \text{ and } B \not\xrightarrow{a})$  or  $(A \not\xrightarrow{a} \text{ and } B \xrightarrow{a})$  then return FALSE;
    return TRUE;

```

Fig. 2. Algorithm EQUIVALENCE(X, Y, G).

By doing that, the algorithm is checking locally for the proof of the nonequivalence of A and B . If the nonequivalence is not discovered and there is nothing to process, i.e., Q is empty, the algorithm returns the value TRUE, meaning $L(A) = L(B)$.

The algorithm EQUIVALENCE is presented in Fig. 2. For technical reasons (to simplify the description of the algorithm) we assume that $First-MP(x, y, H)$ gives ordered pairs in the sense that if $First-MP(x, y, H) = (A, B)$ then $A > B$. For $\alpha \in N^+$ and $a \in \Sigma$, by $\alpha \xrightarrow{a}$ we denote that there is a $\beta \in N^*$ such that $\alpha \rightarrow a\beta$, and by $\alpha \not\xrightarrow{a}$ that there is not. We write $(\alpha, \beta) \xrightarrow{a} (\alpha', \beta')$ to say that $\alpha \rightarrow a\alpha'$ and $\beta \rightarrow a\beta'$.

We present two examples illustrating the execution of algorithm EQUIVALENCE. In the first example the algorithm will determine the equivalence of the two input nonterminals.

Example 15. Consider the execution of EQUIVALENCE(V, Z, G), where $G = (\{a, b\}, \{X, Y, Z, V\}, \{X \rightarrow a, X \rightarrow bXY, Y \rightarrow aX, Y \rightarrow bYY, Z \rightarrow aXX, Z \rightarrow bYXY, V \rightarrow aY, V \rightarrow bYZ\})$.

In the preprocessing phase of the algorithm we compute the lengths of shortest words derivable from the nonterminals, which are $\|X\| = 1$, $\|Y\| = 2$, $\|Z\| = \|V\| = 3$, and then we fix a compatible linear order on the nonterminals, e.g., to $X < Y < Z < V$.

The initialization part sets Q to $\{(V, Z)\}$ and the morphism H to the identity function, i.e., $H = \{X \mapsto X, Y \mapsto Y, Z \mapsto Z, V \mapsto V\}$.

We enter the *while*-loop by choosing an element from Q , i.e., $(\beta_1, \beta_2) = (V, Z)$, and calling $First-MP(\beta_1, \beta_2, H)$. For the presentation only, we explicitly write the values of $H^*(\beta_1)$ and $H^*(\beta_2)$.

In our case $First-MP(\beta_1, \beta_2, H) = (V, Z)$, since $H^*(\beta_1) = V$ and $H^*(\beta_2) = Z$.

The next step is to compute the quotient of the first symbol V by the second symbol Z . Since $\|V\| = \|Z\| = 3$, $quot(V, Z) = \varepsilon$.

Using the values computed by $First-MP()$ and then by $quot()$, the algorithm updates the morphism H by setting $H[V] = Z\varepsilon = Z$.

In the last step of the *while*-loop, the algorithm computes derivations over the same terminal symbols for both V and its new decomposition value Z , i.e., $(V, Z) \xrightarrow{a} (Y, XX)$ and $(V, Z) \xrightarrow{b} (YZ, YXY)$. The resulting pairs (Y, XX) and (YZ, YXY) are added to Q .

The complete trace of the execution of the algorithm is summarized in Fig. 3. The underlined elements in set Q are the pairs (β_1, β_2) considered by the algorithm during the iteration.

Values at the beginning of the iteration					Function calls			
Q	$H(X)$	$H(Y)$	$H(Z)$	$H(V)$	$H^*(\beta_1)$	$H^*(\beta_2)$	<i>First-MP</i>	<i>quot</i>
(V, Z)	X	Y	Z	V	V	Z	(V, Z)	ϵ
$(V, Z), (Y, XX),$ (YZ, YXY)	X	Y	Z	Z	Z	Z	<i>nil</i>	—
$(Y, XX), (YZ, YXY)$	X	Y	Z	Z	Y	XX	(Y, X)	X
$(Y, XX), (YZ, YXY),$ $(X, X), (YY, XYY)$	X	XX	Z	Z	XX	XX	<i>nil</i>	—
$(YZ, YXY), (X, X),$ (YY, XYY)	X	XX	Z	Z	XXZ	$XXXXX$	(Z, X)	XX
$(YZ, YXY), (X, X),$ $(YY, XYY),$ $(XX, XX),$ $(YXY, XYYXX)$	X	XX	XXX	Z	$XXXXX$	$XXXXX$	<i>nil</i>	—
$(X, X), (YY, XYY),$ $(XX, XX),$ $(YXY, XYYXX)$	X	XX	XXX	Z	X	X	<i>nil</i>	—
$(YY, XYY),$ $(XX, XX),$ $(YXY, XYYXX)$	X	XX	XXX	Z	$XXXX$	$XXXX$	<i>nil</i>	—
$(XX, XX),$ $(YXY, XYYXX)$	X	XX	XXX	Z	XX	XX	<i>nil</i>	—
$(YXY, XYYXX)$	X	XX	XXX	Z	$XXXXX$	$XXXXX$	<i>nil</i>	—
empty	X	XX	XXX	Z				

Fig. 3. Execution of EQUIVALENCE(V, Z, G) proving $L_G(V) = L_G(Z)$.

Values at the beginning of the iteration					Function calls			
Q	$H(X)$	$H(Y)$	$H(Z)$	$H(V)$	$H^*(\beta_1)$	$H^*(\beta_2)$	<i>First-MP</i>	<i>quot</i>
(V, Z)	X	Y	Z	V	V	Z	(V, Z)	ϵ
$(V, Z), (XX, Y),$ (ZV, XYZ)	X	Y	Z	Z	Z	Z	<i>nil</i>	—
$(XX, Y), (ZV, XYZ)$	X	Y	Z	Z	XX	Y	(Y, X)	X
$(XX, Y), (ZV, XYZ),$ $(X, X), (XYX, ZXX)$	X	XX	Z	Z	XX	XX	<i>nil</i>	—
$(ZV, XYZ), (X, X),$ (XYX, ZXX)	X	XX	Z	Z	ZZ	$XXXXZ$	(Z, X)	Y
$(ZV, XYZ), (X, X),$ $(XYX, ZXX), (Y, Y),$ (XYZ, ZXY)	X	XX	XY	Z	$XXXXXX$	$XXXXXX$	<i>nil</i>	—
$(X, X), (XYX, ZXX),$ $(Y, Y), (XYZ, ZXY)$	X	XX	XY	Z	X	X	<i>nil</i>	—
$(XYX, ZXX), (Y, Y),$ (XYZ, ZXY)	X	XX	XY	Z	$XXXX$	$XXXX$	<i>failure</i>	

Fig. 4. Execution of EQUIVALENCE(V, Z, G) proving $L_G(V) \neq L_G(Z)$.

After 10 iterations of the *while*-loop, the set Q becomes empty, therefore the two nonterminals V and Z are equivalent.

We now present an example of the execution of algorithm EQUIVALENCE resulting in failure.

Example 16. Let us consider the execution of EQUIVALENCE(V, Z, G), where $G = (\{a, b\}, \{X, Y, Z, V\}, \{X \rightarrow a, X \rightarrow bZX, Y \rightarrow aX, Y \rightarrow bXYX, Z \rightarrow aY, Z \rightarrow bXYZ, V \rightarrow aXX, V \rightarrow bZV\})$.

The trace of the execution of the algorithm is summarized in Fig. 4.

The last call to *First-MP*(XYX, ZXX, H) has returned *failure*, therefore the two nonterminals V and Z are not equivalent.

Lemma 17. *The algorithm is correct. The algorithm makes $O(n)$ iterations.*

Proof. In each iteration, either a pair of strings is removed from Q , or a nonterminal is “eliminated” and no more than $|\Sigma|$ pairs of strings are inserted into Q . The crucial property is that whenever $H(A) = B\gamma$, then the nonterminals in $B\gamma$ are of smaller rank than A , ensuring that H is acyclic. Note also that *First-MP* returns a pair (A, B) such that $H(A) = A$, therefore a nonterminal can only be “eliminated” once. After at most $n - 1$ eliminations, *First-MP* will either find that $H^*(\beta_1) = H^*(\beta_2)$ and remove the pair from Q or return *failure*. Thus, the maximum number of iterations is $O(n)$.

Correctness follows from Lemma 10. \square

Corollary 18. *The algorithm $\text{EQUIVALENCE}(X, Y, G)$ works in time $O(nF(n))$, where n is the size of G , and $F(n)$ is the complexity of the first mismatch-pair problem.*

Lemma 19. *Every instance of $\text{First-MP}(\alpha, \beta, H)$ in $\text{EQUIVALENCE}(X, Y, G)$ can be solved in time:*

- (1) $O(n^6 \log^2 n)$ and
- (2) $O(n^4 \text{polylog } v(G))$,

where n is the size of G , and $v(G)$ is the length of a shortest string derivable from a nonterminal, maximized over all nonterminals.

Proof. In the proof we use twice Lemma 11.

1. Assume H is an acyclic morphism over N , where $n = |N|$ such that $|H(A)| \leq n$ for each A . Then we can construct a morphism H_b such that $H_b^* = H^*$, over a set of $k \leq n^2$ nonterminals and with depth $h = O(n \log n)$.

The transformation of the morphism can be done similarly to a balanced transformation into a Chomsky normal form. If $H(A) = B_1 B_2 \dots B_n$ then we introduce $n - 2$ new auxiliary nonterminals to change it into a balanced binary tree generating $B_1 B_2 \dots B_n$ from A . We need $O(n)$ new nonterminals per each original one, altogether the number of nonterminals increases to $O(n^2)$, i.e., k is in $O(n^2)$. However, the depth is changed only logarithmically. Observe that on each top down path in generation we have at most n original nonterminals, all of them should be different, and at most $O(n \log n)$ auxiliary nonterminals, i.e., h is in $O(n \log n)$. Now, point 1 follows from Lemma 11.

2. We can use the technique from [14] which transforms each grammar generating a single word u into a grammar of depth $O(\log |u|)$ by introducing $O(n \text{polylog } n)$ new nonterminals. Then Lemma 11 can be applied. \square

The series of lemmas give directly the following theorem, due to the fact that after binarization of the morphism the number of variables grows at most quadratically and the depth increases only by a logarithmic factor.

Theorem 20. *The algorithm $\text{EQUIVALENCE}(X, Y, G)$ deciding on the equivalence of two nonterminals X and Y in a simple grammar G , works in time $O(n^7 \log^2 n)$ and $O(n^5 \text{polylog } v(G))$, where n is the size of G , and $v(G)$ is the length of a shortest string derivable from a nonterminal, maximized over all nonterminals.*

6. Conclusion

We have given an efficient algorithm converting any simple grammar to its canonical representation called prime normal form (PNF).

The second result of the paper is the improvement of the complexity of the best existing algorithm verifying the equivalence of simple languages. The worst-case complexity of this algorithm remains quite high. Nevertheless, the approach is applicable in practice in the context of network packet filtering and classification.

The simple language equivalence procedure is the most expensive step of the algorithm finding the PNF representation of simple grammars. Despite the improvement of this procedure, proposed in our paper, the PNF algorithm works in exponential time in the worst case, since its output may be of exponential size. However, this theoretical limitation does not seem to occur in practice.

One interesting open problem is to propose a canonical representation of a simple grammar, and an algorithm computing it, such that the size of this representation is polynomial in the size of the original grammar. However, this would require keeping the internal data structures of such algorithm in a compressed form.

Acknowledgment

The research of the first three authors was supported by NSERC and the research of the fourth author was supported by the Grants KBN 4T11C04425 and ITR-CCR-0313219.

References

- [1] D. Caucal, A fast algorithm to decide on simple grammars equivalence, *Optimal Algorithms*, Lecture Notes in Computer Science, Vol. 401, Springer, Berlin, 1989, pp. 66–85.
- [2] B. Courcelle, Une forme canonique pour les grammaires simples deterministes, *RAIRO Inform.* (1974) 19–36.
- [3] B. Courcelle, An axiomatic approach to the Korenjak–Hopcroft algorithms, *Math. Systems Theory* 16 (3) (1983) 191–231.
- [4] J. Czyzowicz, W. Fraczak, A. Pelc, W. Rytter, Prime decompositions of regular prefix codes, *Implementation and Application of Automata, CIAA 2002*, Lecture Notes in Computer Science, Vol. 2608, Springer, Tours, France, 2003, pp. 85–94.
- [5] W. Debski, W. Fraczak, Concatenation state machines and simple functions, *Implementation and Application of Automata, CIAA 2004*, Lecture Notes in Computer Science, Vol. 3317, Springer, Berlin, 2004, pp. 113–124.
- [6] W. Fraczak, A. Podolak, A characterization of s-languages, *Inform. Process. Lett.* 89 (2) (2004) 65–70.
- [7] L. Georgiadis, R.E. Tarjan, Finding dominators revisited: extended abstract, in: *Proc. 15th Annu. ACM–SIAM Symp. on Discrete Algorithms, SODA 2004*, SIAM, Philadelphia, PA, 2004, pp. 869–878.
- [8] M.A. Harrison, *Introduction to Formal Language Theory*, Addison-Wesley, Reading, MA, 1978.
- [9] Y. Hirshfeld, M. Jerrum, F. Moller, A polynomial algorithm for deciding bisimilarity of normed context-free processes, *Theoret. Comput. Sci.* 158 (1–2) (1996) 143–159.
- [10] M. Karpinski, W. Rytter, A. Shinohara, Pattern-matching for strings with short descriptions, in: Z. Galil, E. Ukkonen (Eds.), *Proc. Sixth Annu. Symp. on Combinatorial Pattern Matching*, Vol. 937, Espoo, Finland, Springer, Berlin, 1995, pp. 205–214.
- [11] A.J. Korenjak, J.E. Hopcroft, Simple deterministic languages, in: *Proc. IEEE Seventh Annu. Symp. on Switching and Automata Theory, IEEE Symp. on Foundations of Computer Science*, 1966, pp. 36–46.
- [12] M. Miyazaki, A. Shinohara, M. Takeda, An improved pattern matching for strings in terms of straight-line programs, *J. Discrete Algorithms* 1 (1) (2000) 187–204.
- [13] W. Plandowski, Testing equivalence of morphisms on context-free languages, in: D. Dummy (Ed.), *ESA, Lecture Notes in Computer Science*, Vol. 855, Springer, Berlin, 1994, pp. 460–470.
- [14] W. Rytter, Application of Lempel–Ziv factorization to the approximation of grammar-based compression, *Theoret. Comput. Sci.* 302 (1–3) (2003) 211–222.