

Faster Algorithm for Bisimulation Equivalence of Normed Context-Free Processes

Sławomir Lasota* and Wojciech Rytter**

Institute of Informatics, Warsaw University, Warsaw, Poland

Abstract. The fastest known algorithm for checking bisimulation equivalence of normed context-free processes worked in $O(n^{13})$ time. We give an alternative algorithm working in $O(n^8 \text{polylog } n)$ time. As a side effect we improve the best known upper bound for testing equivalence of simple context-free grammars from $O(n^7 \text{polylog } n)$ to $O(n^6 \text{polylog } n)$.

1 Introduction

Equivalence checking, that is determining whether two systems are equal under a given notion of equivalence, is an important verification problem with a long history. In this paper we consider systems described by context-free grammars. It is well known that language equivalence is undecidable in this class [1]. A decidability result was obtained by Korenjak and Hopcroft [12] for a restricted class of deterministic context-free grammars (*simple grammars*). Remarkably, the language containment is undecidable even for simple grammars [6].

In the context of process algebras, a grammar may be considered as a description of a transition graph rather than a language. The adequate concept of equivalence is then *bisimilarity* (bisimulation equivalence), a notion strictly finer than language equivalence. For graphs generated by context-free grammars, called *context-free processes*, bisimilarity is known to be decidable due to the result of [5]. It has also been demonstrated that bisimilarity is the only equivalence in van Glabbeek's spectrum [7] which is decidable for context-free processes. This places bisimilarity in a very favourable position.

Historically the first decision procedure for bisimilarity on infinite-state systems was given by [3] for a class of *normed* context-free processes, those defined by grammars in which, roughly, each nonterminal generates at least one word. Clearly, language equivalence is still undecidable in this class, as normedness assumption does not facilitate testing language equality. As language equivalence and bisimilarity coincide on deterministic graphs (in fact, the whole van Glabbeek's spectrum collapses), the result of [3] was a strict extension of [12].

Later, decidability was extended to all context-free processes [5].

In the case of normed context-free processes, a number of simplifications of the proof of [3] appeared [4,10], relying on particular decomposition properties of

* Partially supported by the Polish KBN grant No. 4 T11C 042 25 and by European Community Research Training Network GAMES.

** Supported by the Polish KBN grant No. 4 T11C 044 25.

bisimilarity, and yielding an exponential upper bound for bisimilarity checking. A side effect was an improvement for the equivalence of simple grammars, compared to the complexity of the algorithm of [12] which was $\mathcal{O}(n^v)$, where n is the length of the grammar and v is the length of the shortest word generated, in general exponential in n . Independently, Caucal [4] proposed an algorithm for equivalence of simple grammars working in time $\mathcal{O}(n^3v)$.

Huynh and Tian [11] did a next step and proved that complexity of bisimilarity is in NP^{NP} , the second level of the polynomial hierarchy. A first polynomial-time procedure was finally presented by Hirshfeld, Jerrum and Moller [9].

The algorithm in [9] works in time $\mathcal{O}(n^{13})$ and is hence not satisfactory from the practical point of view. This motivated a further research: in [2] an $\mathcal{O}(n^7 \text{polylog } n)$ time algorithm was proposed for the equivalence of simple grammars. In this paper we report a further progress: we give an $\mathcal{O}(n^8 \text{polylog } n)$ time algorithm for bisimilarity on normed context-free processes, thus improving previous $\mathcal{O}(n^{13})$ time of [9]. We believe that our algorithm is conceptually simpler than that of [9]. It is based on the following two insights. First, we avoid an iterative computation of bisimilarity, by a chain of approximants of the greatest fixed point; instead, we are able to reduce the problem of computing the greatest bisimulation to the problem of finding the greatest solution of certain system of boolean equations, and use the linear-time procedure to find this solution. Secondly, we contribute to the algorithmic theory of compressed strings: we develop a fast algorithm for an auxiliary problem on strings called the First Mismatch Problem, working in $\mathcal{O}(n^5 \text{polylog } n)$ time. As a direct corollary, the equivalence of simple grammars can be decided in $\mathcal{O}(n^6 \text{polylog } n)$ time, which beats the complexity of the (fastest known) algorithm of [2].

Context-Free Processes and Bisimilarity. Let Σ be a finite alphabet and $\mathbf{V} = \{X_1, \dots, X_m\}$ a finite set of variables. By a *process definition* Δ we mean a finite set of rules of the form: $X \xrightarrow{a} \alpha$, with $a \in \Sigma$ and $\alpha \in \mathbf{V}^*$. Such process definitions are usually called in the literature *Basic Process Algebra*, or *Context-Free Processes*. The explanation of the latter is that each rule can be seen as a production $X \rightarrow a\alpha$ of a context-free grammar in Greibach normal form. Elements of \mathbf{V}^* are called here *processes*; a variable X can be seen as an elementary process.

Δ defines a transition system: its states are processes $\alpha \in \mathbf{V}^*$; and for each $a \in \Sigma$, there is a transition relation containing triples (α, a, β) , where $a \in \Sigma$ and $\alpha, \beta \in \mathbf{V}^*$, written $\alpha \xrightarrow{a} \beta$. The transition relations are defined by a prefix rewriting: $X\beta \xrightarrow{a} \alpha\beta$ whenever Δ contains a rule $X \xrightarrow{a} \alpha$, and $\beta \in \mathbf{V}^*$.

Definition 1. *Given a binary relation R over \mathbf{V}^* , we say that a pair (α, β) of processes satisfies expansion in R if*

- whenever $\alpha \xrightarrow{a} \alpha'$, there exists some β' with $\beta \xrightarrow{a} \beta'$ and $(\alpha', \beta') \in R$; and
- whenever $\beta \xrightarrow{a} \beta'$, there exists some α' with $\alpha \xrightarrow{a} \alpha'$ and $(\alpha', \beta') \in R$.

A binary relation S satisfies expansion in R if each pair $(\alpha, \beta) \in S$ does. A relation R is a bisimulation if it satisfies expansion in itself. We say that α and β are bisimilar, denoted by $\alpha \sim \beta$, if (α, β) belongs to some bisimulation.

Assume that Δ is *normed*, i.e., for each variable $X \in \mathbf{V}^*$ there is a finite sequence $X \xrightarrow{a_1} \alpha_1 \dots \xrightarrow{a_k} \alpha_k = \varepsilon$ leading from X to the empty process ε . By $|X|$ denote the smallest length of such sequence and call it the *norm* of X (intuitively, $|X|$ is the length of the shortest word generated from X).

We consider the following NORMED-BPA-BISIM PROBLEM:

Instance: A normed Δ and $X, Y \in \mathbf{V}$ with $|X| = |Y|$, $X \neq Y$.

Question: Is $X \sim Y$?

A more general problem of checking whether $\alpha \sim \beta$, for any $\alpha, \beta \in \mathbf{V}^*$, can be easily reduced to the above one. We use notation $\tilde{\mathcal{O}}(f(n))$ for $\mathcal{O}(f(n))$ polylog n in the sequel. The size of Δ , denoted by n , is the sum of lengths of all the rules in Δ . Our main result is the following:

Theorem 1. NORMED-BPA-BISIM PROBLEM can be solved in time $\tilde{\mathcal{O}}(n^8)$.

2 Terminology and Tools Used in the Main Algorithm

The normedness assumption implies that each variable has at least one rule in Δ . We extend additively the norm to all processes: $|\varepsilon| := 0$, $|X\alpha| := |X| + |\alpha|$. Bisimilarity preserves norm, as a sequence of transitions $X \xrightarrow{a_1} \dots \xrightarrow{a_k} \varepsilon$ leading to ε must be necessarily matched by a sequence leading to ε as well:

Proposition 1. If $\alpha \sim \beta$ then $|\alpha| = |\beta|$.

Let Σ , \mathbf{V} and Δ be fixed from now on. We assume that the variables in \mathbf{V} are ordered so that $|X_i| \leq |X_j|$ whenever $i < j$. It is easy to show the following:

Proposition 2. Norms of all variables can be computed in time $\tilde{\mathcal{O}}(n)$.

The following fact is easily derived from the unique decomposition property [9].

Lemma 1. If $\alpha\alpha' \sim \beta\beta'$ and $|\alpha| \geq |\beta|$ then for some γ , $\alpha \sim \beta\gamma$ and $\gamma\alpha' \sim \beta'$.

As a direct corollary we get a cancellation property:

Lemma 2. If $\gamma\alpha \sim \gamma\beta$ then $\alpha \sim \beta$.

We will need a notion of *base*, which is a slight adaptation of [9]. Intuitively, it describes ways of decomposing an elementary process into smaller ones:

Definition 2. A **base** is a set B of pairs $(X_j, X_i\gamma)$, at most one for each pair (X_j, X_i) , such that $i < j$, $\gamma \in \mathbf{V}^*$ and $|X_j| = |X_i\gamma|$.

A base is **full** iff whenever $X_j \sim X_i\beta$, for $j > i$, then $(X_j, X_i\gamma) \in B$, for some $\gamma \sim \beta$.

Note that whenever $(X_j, X_i\gamma) \in B$ then necessarily $X_i\gamma \in \{X_1, \dots, X_{j-1}\}^+$. In the sequel we rely on the following lemma proved in [9]:

Lemma 3. A full base B_0 can be constructed, in time $\mathcal{O}(n^3)$, such that the length of γ is $\mathcal{O}(n)$, for each $(X_j, X_i\gamma) \in B_0$. Furthermore, B_0 contains a pair $(X_j, X_i\gamma)$ for each i, j such that $j > i$.

Remark 1. By cancellation, if $X_j \sim X_i\gamma$, then γ is unique up to bisimilarity. Basing on this observation, the construction of B_0 is by inspecting an arbitrarily chosen sequence of $|X_i|$ norm-reducing moves from X_j . Any process obtained at the end of such a sequence is a good candidate for γ .

Remark 2. Note that B_0 , being full, contains all pairs (X_j, X_i) with $X_j \sim X_i$.

Fix the full base B_0 from now on. Pairs $(X_j, X_i\gamma) \in B_0$ we call *decomposition pairs*, or *d-pairs* in short. A d-pair $(X_j, X_i\gamma)$ will be denoted by z_{ji} .

In the sequel we will treat the d-pairs as boolean variables. The basic intuition will be that $z_{ji} = \mathbf{true}$ just in case when $X_j \sim X_i\gamma$ holds. We will also build the positive boolean formulas on top of d-pairs, by boolean connectives \wedge, \vee and symbols $\mathbf{true}, \mathbf{false}$ (no negation). The empty conjunction (disjunction) is allowed as a formula and understood as \mathbf{true} (\mathbf{false} , respectively).

A *valuation* is a mapping v from B_0 to $\{\mathbf{true}, \mathbf{false}\}$. We extend valuations to formulas in the obvious way.

Definition 3. A boolean equation system is a set of equations of the form

$$z_{ji} = \psi_{ji},$$

one for each $z_{ji} \in B_0$, where ψ_{ji} is a positive boolean formula with variables from B_0 . A solution is any valuation v such that $v(z_{ji}) = v(\psi_{ji})$ for each $z_{ji} \in B_0$.

Valuations are in one-to-one correspondence with subsets of B_0 : for $B \subseteq B_0$, a corresponding valuation v_B assigns true to a variable z_{ji} if and only if $z_{ji} \in B$. Each boolean equation system has the greatest solution \bar{B} : start with $B = B_0$ and iteratively update B by removing z_{ji} from B if $v_B(\psi_{ji}) = \mathbf{false}$, until B eventually stabilizes yielding \bar{B} . A relevant observation is a folklore (see e.g. [8]):

Lemma 4. The greatest solution of a boolean equation system can be computed in time linear wrt. the size of the system.

The overall idea underlying the algorithm is as follows. Bisimilarity is the greatest bisimulation, i.e., \sim satisfies expansion in itself. Hence, by Knaster-Tarski fix-point theorem, it is also the greatest fixed point: $\alpha \sim \beta$ if and only if (α, β) satisfies expansion in \sim . One crucial insight is that computing this greatest fixed point can be reduced to finding the greatest solution of certain system of boolean equations. Another insight is that the system of equations can be constructed effectively (due to Lemma 5 below). These two insights allowed us to obtain an algorithm working in time $\tilde{O}(n^8)$.

Consider any sub-base $B \subseteq B_0$ and two processes $\alpha, \beta \in \mathbf{V}^*$ of equal norm. We write $\alpha =_B \beta$ if α and β can be shown equal by applying any number of substitutions $X_j \mapsto X_i\gamma$, where $(X_j, X_i\gamma) \in B$. Formally:

Definition 4. For $B \subseteq B_0$, let $=_B$ be the smallest symmetric relation over processes containing all identical pairs $\alpha =_B \alpha$ and such that if $(X_j, X_i\gamma) \in B$ and $\alpha X_i\gamma =_B \delta$ then $\alpha X_j =_B \delta$.

Example 1. If $|A| > |B| > |C| > |D| > |E|$ and $B_0 = \{(A, BBD), (A, ECEB), (C, DE), (B, CD), (B, DED), \dots\}$ then we have $AEBBBD =_B BBCCDA$ for $B = \{(A, BBD), (B, CD), (C, DE)\}$, due to the derivations of the same string:

$$\begin{aligned}
 AEBBBD &\xrightarrow{A=BBD} BBDEBBBD \xrightarrow{B=CD} BBDECDBBD, \\
 BBCCDA &\xrightarrow{C=DE} BBDECDA \xrightarrow{A=BBD} BBDECDBBD.
 \end{aligned}$$

Note that B is inclusion-minimal, i.e., there is no $B' \subsetneq B$ with $AEBBBD =_{B'} BBCCDA$.

As B_0 contains a pair $(X_j, X_i\gamma)$ for each i and j with $j > i$, it follows that for α, β of equal norm some sub-base B with $\alpha =_B \beta$ always exists; in particular $\alpha =_{B_0} \beta$. The relevant issue will be to find such B possibly small. For future reference, let $B_\sim \subseteq B_0$ be the set of all d-pairs $(X_j, X_i\gamma)$ satisfying $X_j \sim X_i\gamma$.

Definition 5. We say that B is a **matching sub-base** for α, β iff it is an inclusion-minimal subset of B_0 such that (i) $\alpha =_B \beta$ and (ii) $\alpha \sim \beta$ implies $B \subseteq B_\sim$.

In particular, by inclusion-minimality the matching sub-base for α, α is \emptyset . Condition (ii) says that a bisimilar pair can be shown equal by using only bisimilar substitutions. This property will follow by the unique decomposition (cf. Lemma 1) and by our intricate construction of the sub-base in the proof of Lemma 5. The proof of the lemma is postponed to Section 4.

Lemma 5. For any processes α, β of length $O(n)$ with $|\alpha| = |\beta|$, we can compute in $\tilde{O}(n^6)$ time a matching sub-base $B_{\alpha,\beta}$ containing $O(n)$ d-pairs.

3 The Main Algorithm

Basing on $B_{\alpha,\beta}$, we define a *matching formula* for α, β , denoted by $\phi_{\alpha,\beta}$, as follows: if $|\alpha| = |\beta|$, then $\phi_{\alpha,\beta}$ is a conjunction of all d-pairs z_{ji} belonging to $B_{\alpha,\beta}$, otherwise $\phi_{\alpha,\beta}$ is **false**.

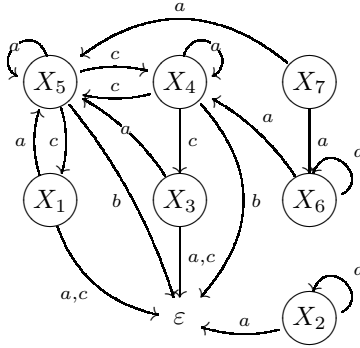
In the algorithm we construct a boolean equation system containing an equation $z_{ji} = \psi_{ji}$ for each z_{ji} and then apply Lemma 4. The intuition is that formula ψ_{ji} expresses the property that the d-pair z_{ji} satisfies expansion in \sim . However, instead of directly referring to $\alpha \sim \beta$ in ψ_{ji} , we will prefer to use formulas $\phi_{\alpha,\beta}$ as subformulas in ψ_{ji} , relying on condition (ii) in Definition 5.

Let $z_{ji} = (X_j, X_i\gamma)$. The formula ψ_{ji} is defined as follows:

$$\psi_{ji} = \bigwedge_a \left(\bigwedge_{\beta} \bigvee_{\alpha} \phi_{\beta,\alpha\gamma} \wedge \bigwedge_{\alpha} \bigvee_{\beta} \phi_{\beta,\alpha\gamma} \right), \tag{1}$$

where a ranges over Σ , α ranges over $\{\delta : X_i \xrightarrow{a} \delta \in \Delta\}$ and β ranges over $\{\delta : X_j \xrightarrow{a} \delta \in \Delta\}$.

Example 2. As an illustration, consider Δ containing variables X_1, \dots, X_7 and the rules $X_5 \xrightarrow{c} X_1, X_7 \xrightarrow{a} X_5$, etc., as shown in the picture:



Δ is essentially finite-state. All variables have norm 1 except $|X_7| = |X_6| = 2$. Let $B_0 = \{z_{76} = (X_7, X_6), z_{72} = (X_7, X_2X_5), z_{65} = (X_6, X_5X_4), z_{64} = (X_6, X_4X_4), z_{62} = (X_6, X_2X_4), z_{54} = (X_5, X_4), z_{32} = (X_3, X_2), \dots\}$. $B_{\sim} = \{(X_3, X_1), (X_5, X_4), (X_7, X_6), (X_6, X_2X_4), (X_7, X_2X_5)\}$. For instance, the equation for z_{76} is:

$$z_{76} = (\phi_{X_5, X_4} \vee \phi_{X_5, X_6}) \wedge (\phi_{X_6, X_4} \vee \phi_{X_6, X_6}) \wedge (\phi_{X_6, X_4} \vee \phi_{X_5, X_4}) \wedge (\phi_{X_6, X_6} \vee \phi_{X_5, X_6}).$$

The first conjunct describes possible matchings for $X_7 \xrightarrow{a} X_5$, by $X_6 \xrightarrow{a} X_4$ or $X_6 \xrightarrow{a} X_6$; the second conjunct describes possible matchings for $X_7 \xrightarrow{a} X_6$, etc. Note for instance that $\phi_{X_5, X_4} = z_{54}$. Furthermore $\phi_{X_6, X_6} = \mathbf{true}$ as $B_{X_6, X_6} = \emptyset$; and $\phi_{X_5, X_6} = \mathbf{false}$ as $|X_5| \neq |X_6|$. After simplification we derive: $z_{76} = z_{54}$. Similarly we derive:

$$z_{54} = z_{54} \wedge (z_{51} \vee z_{31}) \wedge (z_{54} \vee z_{43}) \wedge (z_{31} \vee z_{43}) \wedge (z_{51} \vee z_{54}).$$

The greatest solution of the equation system is $\{z_{31}, z_{54}, z_{76}, z_{62}, z_{72}\}$.

Algorithm Normed-BPA-Bisim(Δ, X, Y); // $|X| = |Y|$

- (1) compute a full base B_0 ; // c.f. Lemma 3
- (2) **for each** $(X_j, X_i\gamma) \in B_0$ and $a \in \Sigma$ **do**
 for each $X_j \xrightarrow{a} \beta$ and $X_i \xrightarrow{a} \alpha$ in Δ **do**
 compute the formula $\phi_{\beta, \alpha\gamma}$;
- (3) **for each** $z_{ji} = (X_j, X_i\gamma) \in B_0$ **do**
 construct the boolean expression (1);

// this yields a boolean equation system S

- (4) compute the greatest solution $\bar{B} \subseteq B_0$ of S ;
- (5) **return** $[(X, Y) \in \bar{B}]$.

Lemma 6. *The algorithm works in time $\tilde{O}(n^8)$.*

Proof. Step (1) requires time $\mathcal{O}(n^3)$, by Lemma 3. The total number of invocations of the procedure computing $\phi_{\beta,\alpha\gamma}$ in step (2) is $\mathcal{O}(n^2)$. Hence, step (2) can be completed in time $\tilde{O}(n^8)$, by Lemma 5, and this is dominating the total cost. The size of the boolean equation system built in step (3) is $\mathcal{O}(n^3)$: indeed, the length of each subformula $\phi_{\beta,\alpha\gamma}$ is $\mathcal{O}(n)$, and there is a quadratic number of such formulas in the right-hand sides of equations. The equation system can be constructed in time $\mathcal{O}(n^3)$ and its greatest solution can be computed in the same time, by Lemma 4. □

Lemma 7. *The algorithm is correct.*

Proof. Correctness follows directly from the equality $B_{\sim} = \bar{B}$, which is shown below in two steps.

We show $B_{\sim} \subseteq \bar{B}$ first. Denote by ψ_{ji} the right-hand side of equation (1) for $z_{ji} = (X_j, X_i\gamma)$. \bar{B} , as the greatest solution of the equation system, is the greatest subset $B \subseteq B_0$ such that $z_{ji} \in B$ (i.e., $v_B(z_{ji}) = \mathbf{true}$) implies $v_B(\psi_{ji}) = \mathbf{true}$. Hence we will only show that $z_{ji} \in B_{\sim}$ (i.e., $X_j \sim X_i\gamma$) implies $v_{B_{\sim}}(\psi_{ji}) = \mathbf{true}$.

Indeed. If $X_j \sim X_i\gamma$ then this pair satisfies expansion in \sim . Consider any pair $\beta \sim \alpha\gamma$ relevant for the expansion, with $X_j \xrightarrow{a} \beta$ and $X_i \xrightarrow{a} \alpha$ for some a . By point (ii) in Definition 5 we know that all d-pairs appearing in the matching formula $\phi_{\beta,\alpha\gamma}$ are bisimilar. As this applies to any pair $\beta \sim \alpha\gamma$, the right-hand side formula ψ_{ji} is true under valuation $v_{B_{\sim}}$, as required.

Now we will prove $\bar{B} \subseteq B_{\sim}$. Consider any solution B of the equation system and any d-pair $z_{ji} = (X_j, X_i\gamma) \in B$. Thus the right-hand side ψ_{ji} evaluates to \mathbf{true} under valuation v_B . Consider any matching formula $\phi_{\beta,\alpha\gamma}$ appearing in ψ_{ji} . If it is true under valuation v_B then $B_{\beta,\alpha\gamma} \subseteq B$ and hence $\beta =_B \alpha\gamma$, by point (i) in Definition 5. Hence, by the very construction of ψ_{ji} on top of the matching formulas $\phi_{\beta,\alpha\gamma}$ it follows that $(X_j, X_i\gamma)$ satisfies expansion in $=_B$.

But $=_B$ is clearly contained in $\overset{B}{\equiv}$, the smallest congruence containing B . As the d-pair z was chosen arbitrarily, we have shown that B satisfies expansion in $\overset{B}{\equiv}$, i.e., B is a so called *Caucal base*, or *self-bisimulation* [4]. By a standard argument $\overset{B}{\equiv} \subseteq \sim$, and hence $B \subseteq B_{\sim}$. As all this was said for an arbitrary solution, in particular $\bar{B} \subseteq B_{\sim}$ as required. □

4 Computing a Matching Sub-base

This section contains the construction of sub-bases $B_{\alpha,\beta}$ and the proof of Lemma 5. To this aim we need to refine the concept of base a little bit further.

Definition 6. *A production is any pair (X_i, γ) , written $X_i \rightarrow \gamma$, such that $|X_i| = |\gamma|$ and $\gamma \in \{X_1, \dots, X_{i-1}\}^+$. A **decomposition grammar** (*d-grammar*) G is any set of productions $X_i \rightarrow \gamma$, at most one for each variable X_i .*

Let $\mathbf{V}(G)$ denote the set of all X_i such that $(X_i \rightarrow \gamma) \in G$ for some γ . Variables in $\mathbf{V}(G)$ and $\mathbf{V} \setminus \mathbf{V}(G)$ are called *nonterminal* and *terminal symbols*, respectively.

Each nonterminal X_i has precisely one production, hence generates a single nonempty word $G(X_i) \in (\mathbf{V} \setminus \mathbf{V}(G))^+$. We extend this to words $\alpha \in \mathbf{V}^*$ in the obvious way and write $G(\alpha)$ for the single word in $(\mathbf{V} \setminus \mathbf{V}(G))^*$ generated from α . A d-grammar G induces an equivalence $=_G$ over \mathbf{V}^* : $\alpha =_G \beta$ iff α and β generate the same word, i.e., $G(\alpha) = G(\beta)$.

Assume $|\alpha| = |\beta|$ hence $|G(\alpha)| = |G(\beta)|$. One of $G(\alpha)$, $G(\beta)$ can not be a proper prefix of the other. Hence, if $\alpha \neq_G \beta$, the words $G(\alpha)$ and $G(\beta)$ must have the left-most mismatching pair of variables. Formally: there is some i, j, γ such that $i \neq j$, γX_i is a prefix of $G(\alpha)$ and γX_j is a prefix of $G(\beta)$. W.l.o.g. assume $i < j$. The pair (X_j, X_i) is called the *first mismatch-pair of α and β wrt. G* and denoted by $\text{First-MP}(\alpha, \beta, G)$. If $|\alpha| = |\beta|$ and $\alpha =_G \beta$, we put $\text{First-MP}(\alpha, \beta, G) = \text{nil}$. We define **FIRST MISMATCH PROBLEM**:

Input: A d-grammar G and two processes $\alpha, \beta \in \mathbf{V}^*$ of equal norm.
Output: $\text{First-MP}(\alpha, \beta, G)$.

Lemma 8. **FIRST MISMATCH PROBLEM** can be solved in time $\tilde{\mathcal{O}}(n^5)$, if the lengths of α , β and all productions $X_i \rightarrow \gamma$ in G are in $\mathcal{O}(n)$.

Section 5 is devoted to the proof of this lemma. In the rest of this section Lemma 8 will be used to prove Lemma 5.

```

Computation of  $B_{\alpha,\beta}$ : //  $|\alpha| = |\beta|$ 
 $G := \emptyset$ ;
while  $\text{First-MP}(\alpha, \beta, G) \neq \text{nil}$  do
     $(X_j, X_i) := \text{First-MP}(\alpha, \beta, G)$ ;
    // let  $\gamma$  be the unique process such that  $(X_j, X_i\gamma) \in B_0$ 
     $G := G \cup \{X_j \rightarrow X_i\gamma\}$ ;
 $B_{\alpha,\beta} := G$ .
    
```

Note that G is always a d-grammar in the course of the computation, as whenever a production $X_j \rightarrow X_i\gamma$ is added to G , X_j has no production yet in G .

For instance, for B_0 and the two processes considered in Example 1 we obtain $B_{AEBBBB, BBCCDA} = \{(A, BBD), (C, DE), (B, DED)\}$.

Proof (of Lemma 5). The computation of $B_{\alpha,\beta}$ needs $\mathcal{O}(n)$ calls to **FIRST MISMATCH PROBLEM**, hence it completes in $\tilde{\mathcal{O}}(n^6)$ time.

Now we will show that $B_{\alpha,\beta}$ is a matching sub-base for α, β , in the sense of Definition 5. By the very construction, $B_{\alpha,\beta}$ is an inclusion-minimal d-grammar with $\alpha =_{B_{\alpha,\beta}} \beta$. Furthermore, clearly $B_{\alpha,\beta} \subseteq B_0$. Hence, the equivalence $=_{B_{\alpha,\beta}}$ induced by $B_{\alpha,\beta}$ as a d-grammar is a special case of the relation $=_B$ induced by a sub-base $B \subseteq B_0$, cf. Definition 4. This completes the proof of condition (i) in Definition 5.

For condition (ii), assume $\alpha \sim \beta$. We will show that $B_{\alpha,\beta} \subseteq B_{\sim}$. It is sufficient to prove that each production $X_j \rightarrow X_i\gamma$ added to G in the course of the computation satisfies $X_j \sim X_i\gamma$.

Assume therefore $G \subseteq B_{\sim}$ (this implies $G(\alpha) \sim G(\beta)$) and $G(\alpha) \neq G(\beta)$. We need to consider the first mismatch-pair of α and β wrt. G , say (X_j, X_i) . By Lemma 2 we can ignore the matching prefixes of $G(\alpha)$ and $G(\beta)$, and then by Lemma 1 applied to $\alpha = X_j$ and $\beta = X_i$, we conclude that for some γ' it holds $X_j \sim X_i\gamma'$. Let γ be the unique process for which $(X_j, X_i\gamma) \in B_0$. As B_0 is full, it satisfies $\gamma \sim \gamma'$. As a consequence, $X_j \sim X_i\gamma$ as required. \square

5 Algorithm for the First Mismatch Problem

Let G be a given d-grammar. We start with the problem of **equality-testing**: for two nonterminals S_1, S_2 test if they generate the same string. We identify informally the names of nonterminals with their values, so it can be written as $S_1 = S_2$ instead of $S_1 =_G S_2$.

If $A \rightarrow A_1A_2 \dots A_r$, then by cut-points, or decomposition points, we mean the positions $|A_1|, |A_1| + |A_2|, \dots, |A_1| + |A_2| + |A_3| + \dots + |A_{r-1}|$, see Figure 1.

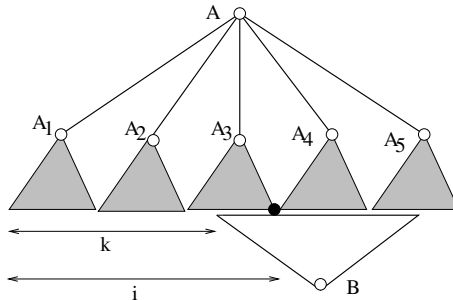


Fig. 1. Assume there is a production $A \rightarrow A_1A_2 \dots$. In this case $i = |A_1| + |A_2| + |A_3|$ is the third cut-point of A (black circle) and k is the distance between the possible occurrence of B and the beginning of A . Validity of the overlap item (A, B, i, k) is equivalent to $B = A[k + 1 \dots k + |B|]$.

An **overlap item** is a 4-tuple (A, B, i, k) such that i is a decomposition point of A and k is a beginning position of a *potential* occurrence of B in A which overlaps i , see Figure 1. Overlapping means that the occurrence of B is *touching* the point i , i.e., $k \leq i \leq k + |B| \leq |A|$. This overlap item is said to be **valid** iff $B = A[k + 1 \dots k + |B|]$.

The equality of two nonterminals S_1, S_2 is equivalent to the overlap item $\alpha_0 = (S_1, S_2, FirstDecPoint(S_1), 0)$, where $FirstDecPoint(S_1)$ denotes the first decomposition point of S_1 . Let us fix in this section S_1, S_2 and α_0 .

We say that a set of items $\Gamma = \{\gamma_1, \gamma_2, \dots, \gamma_p\}$ covers an item β iff

$$[\alpha_0 \implies (\gamma_1 \ \& \ \gamma_2 \ \& \ \dots \ \& \ \gamma_p)] \quad \text{and} \quad [(\gamma_1 \ \& \ \gamma_2 \ \& \ \dots \ \& \ \gamma_p) \implies \beta].$$

Observe that if T covers α_0 then equality of $S_1 = S_2$ is equivalent to validity of all items in T .

Recall that nonterminals are ordered with respect to the increasing norms of their values. If \mathcal{D} is a set of items then *DeleteLexMax*(\mathcal{D}) returns lexicographically maximal element of \mathcal{D} and removes it from \mathcal{D} . An item is **atomic** iff the nonterminals occurring in this item generate only terminal symbols. We can test validity of each individual atomic item in constant time.

We describe the basic functions in the equality testing.

The function *SubtleInsert*(β, \mathcal{D}) inserts β into \mathcal{D} in $\tilde{O}(n)$ time. For every nonterminal B and every cut-point of A we keep only at most three occurrences of B overlapping A on this cut-point. Correctness follows from the fact that the set of occurrence of the same string overlapping a given cut-point is a single arithmetic progression. The function *SubtleInsert* inserts only if it is necessary, and if it inserts β and there are already three occurrences overlapping the same cut-point then one of them is removed.

Next we describe how to implement the function GENERATE. Let α be a non-atomic item. *GENERATE*(α) is a set of items satisfying the following property: **(1)** it is of size $O(n)$; **(2)** it contains only items lexicographically smaller than α ; **(3)** it covers α .

```

function EqTest( $S_1, S_2$ ); //  $|S_1| = |S_2|$ 

  for each production do
    sort the set of its cut-points;
   $\alpha_0 := (S_1, S_2, FirstDecPoint(S_1), 0)$ ;  $\mathcal{D} := \{\alpha_0\}$ ;

  while  $\mathcal{D}$  contains a non-atomic item do
     $\alpha := DeleteLexMax(\mathcal{D})$ ;
    for each  $\beta \in GENERATE(\alpha)$  SubtleInsert( $\beta, \mathcal{D}$ );

  Comment:  $\mathcal{D}$  covers  $\alpha_0$  and consists only of atomic items;

  return  $[(\forall \alpha \in \mathcal{D}) valid(\alpha)]$ 
    
```

Lemma 9. Let $\alpha = (A, B, i, k)$. We can compute *GENERATE*(α) satisfying the conditions (1-3) above in $\tilde{O}(n)$ time.

Proof. In the proof we use temporarily other type of items: an *internal item* is a triple (X, Y, t) , where t is a potential occurrence of B in A , not necessarily overlapping a cut-point of A . Assume $\alpha = (A, B, i, k)$, and there is a production $B \rightarrow B_1 B_2 \dots B_r$. We can locate each of B_i in A and we have a set of internal items (A, B_1, k) , $(A, B_2, k + i_1)$, \dots , $(A, B_r, k + i_1 + \dots + i_{r-1})$. We can design a subroutine *overlapify*($A, B_p, k + i_1 + \dots + i_{p-1}$), this subroutine finds an overlap item which covers this internal item. It is possible to design such a subroutine *overlapify* which, applied to all internal items, finds together

in $\tilde{O}(n)$ time the set of overlap items covering them. Consequently it finds a set of smaller overlap items covering (A, B, i, k) . The leftmost and rightmost internal item is *overlapped* in $\tilde{O}(n)$ time, all others are processed in logarithmic time, using the sorted order of cut-points and a kind of binary search. This set of overlap items is returned by the function GENERATE. Figure 2 shows how the overlap item $(A, B, *, *)$ is decomposed into the set of internal items $(A, B1, *)$, $(A, B2, *)$, \dots $(A, B6, *)$. The leftmost and rightmost internal items are covered by finding the lowest common ancestors (denoted by X, U in Figure 2) of the endpoints of $B1$ and $B6$. This takes $O(n)$ time. All other internal items are covered in logarithmic time, after merging cut-points of B and the set \mathcal{S} (illustrated as small darkened circles in Figure 2) of decomposition points of nonterminals branching from the paths from the root to the nodes X, U . It is crucial that after sorting cut-points for each nonterminal we can preprocess \mathcal{S} in such a way that binary searching (because \mathcal{S} will be sorted) can be done in logarithmic time. The set \mathcal{S} is of size $O(n^2)$ but it consists of $O(n)$ groups of cut-points of nonterminals on the branches from A to X or U . Each group is sorted. Also the beginning and ending positions ($O(n)$ together) of these groups can be first sorted. We omit the details.

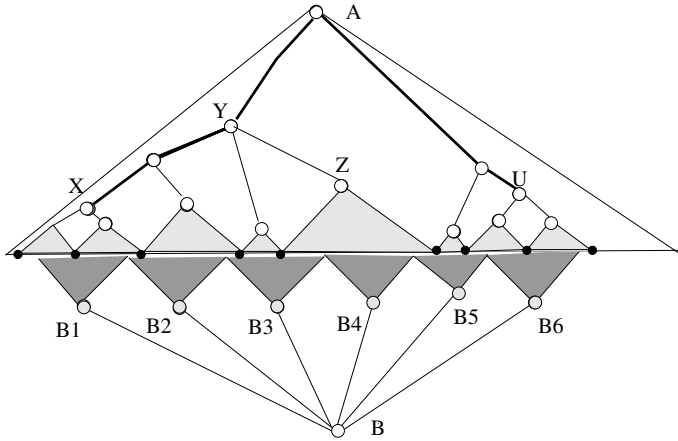


Fig. 2. $(A, B1, *)$ is covered by $(X, B1, *, *)$, $(A, B4, *)$ is covered by $(B, Z, *, *)$ and $(A, B3, *)$ is covered by $(Y, B3, *, *)$, the *'s denote corresponding positions in the items. The set \mathcal{S} consists of small darkened circles.

Proof (of Lemma 8). First we analyse the function GENERATE. The total number of overlap items is $O(n^3)$, we process each item only once with the function GENERATE, it takes $\tilde{O}(n)$ time per single item, according to Lemma 9. Altogether the complexity of equality test is $\tilde{O}(n^4)$.

Now we can find the first mismatch using the algorithm for equality testing and a kind of binary search. We need at most $O(n)$ instances of equality testing, since the depth of the grammar is $O(n)$. Hence the overall time is $\tilde{O}(n^5)$.

6 Equivalence of Simple Grammars

The class of simple grammar is the largest class of context-free grammars for which equivalence problem can be tested in deterministic polynomial time. This is nontrivial since inclusion problem for this class of grammars is undecidable. A simple grammar is a context-free grammar in Greibach normal form, such that whenever $A \rightarrow a \alpha$ and $A \rightarrow a \beta$ then $\alpha = \beta$. The main component in the $\tilde{O}(n^7)$ time algorithm of [2] is the computation of the first mismatch problem in $\tilde{O}(n^6)$ time. As we improved this to $\tilde{O}(n^5)$ we have immediately the following result.

Theorem 2. *Equivalence of simple grammars can be tested in $\tilde{O}(n^6)$ time.*

References

1. Y. Bar-Hillel, M. Perles, and S. Shamir. On formal properties of simple phrase structure grammars. *Zeitschrift fuer Phonetik, Sprachwissenschaft, und Kommunikationsforschung*, 14:143–177, 1961.
2. C. Bastien, J. Czyżowicz, W. Fraczak, and W. Rytter. Prime normal form and equivalence of simple grammars. In *Proc. CIAA '05*, volume 3845 of *LNCS*, pages 79–90. Springer-Verlag, 2005.
3. J. Beaten, J. Bergstra, and J. Klop. Decidability of bisimulation equivalence for processes generating context-free languages. In *Proc. PARLE'87*, volume 259 of *LNCS*, pages 94–113. Springer-Verlag, 1987.
4. D. Caucal. Graphes canoniques des graphes algébriques. *Informatique Théorique et Applications (RAIRO)*, 24(4):339–352, 1990.
5. S. Christensen, Y. Hirshfeld, and C. Stirling. Bisimulation equivalence is decidable for all context-free processes. *Information and Computation*, 12(2):143–148, 1995.
6. E.P. Friedman. The inclusion problem for simple languages. *Theoretical Computer Science*, 1:297–316, 1976.
7. R. v. Glabbeek. The linear time - branching time spectrum. In *Proc. CONCUR'90*, pages 278–297, 1990.
8. J. Groote and M. Keinänen. A Sub-quadratic Algorithm for Conjunctive and Disjunctive BESs. CS-Report 04-13, Department of Mathematics and Computer Science, Technische Universiteit Eindhoven, June 2004.
9. Y. Hirshfeld, M. Jerrum, and F. Moller. A polynomial algorithm for deciding bisimilarity on normed context-free processes. *Theoretical Computer Science*, 15:143–159, 1996.
10. H. Huettel and C. Stirling. Actions speak louder than words: proving bisimilarity for context-free processes. In *Proc. LICS'91*, pages 376–386. IEEE Computer Society Press, 1991.
11. D. Huynh and L. Tian. Deciding bisimilarity of normed context-free processes is in Σ_2^P . *Theoretical Computer Science*, 123:183–197, 1994.
12. A. Korenjak and J. Hopcroft. Simple deterministic languages. In *Proc. 7th Annual IEEE Symposium on Switching and Automata Theory*, pages 36–46, 1966.
13. W. Plandowski. Testing equivalence of morphisms on context-free languages. In *Proc. ESA'94*, volume 855 of *LNCS*, pages 460–470. Springer-Verlag, 1994.
14. A. Shinohara, M. Miyazaki, and M. Takeda. An improved pattern-matching for strings in terms of straight-line programs. *Journal of Discrete Algorithms*, 1(1):187–204, 2000.