



HAL
open science

Right-Universality of Visibly Pushdown Automata

Véronique Bruyère, Marc Ducobu, Olivier Gauwin

► **To cite this version:**

Véronique Bruyère, Marc Ducobu, Olivier Gauwin. Right-Universality of Visibly Pushdown Automata. 4th International Conference on Runtime Verification, Sep 2013, Rennes, France. pp.76-93, 10.1007/978-3-642-40787-1_5 . hal-00946193

HAL Id: hal-00946193

<https://hal.science/hal-00946193v1>

Submitted on 13 Feb 2014

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Right-Universality of Visibly Pushdown Automata

Véronique Bruyère¹, Marc Ducobu¹, and Olivier Gauwin²

¹ University of Mons, UMONS, Belgium

² University of Bordeaux, LaBRI, France

Abstract. Visibly pushdown automata (VPAs) express properties on structures with a nesting relation such as program traces with nested method calls. In the context of runtime verification, we are interested in the following problem: given u , the beginning of a program trace, and \mathcal{A} , a VPA expressing a property to be checked on this trace, can we ensure that any extension uv of u will be accepted by \mathcal{A} ? We call this property right-universality w.r.t. u . We propose an online algorithm detecting at the earliest position of the trace, whether this trace is accepted by \mathcal{A} . The decision problem associated with right-universality is ExpTime-complete. Our algorithm uses antichains and other optimizations, in order to avoid the exponential blow-up in most cases. This is confirmed by promising experiments conducted on a prototype implementation.

1 Introduction

Program traces describe the control flows of program executions, including subroutine calls and returns. Some properties of program traces are specific to the nesting structure of calls and returns, as for instance reentrant locks [BYBC10]. These properties are not captured by regular languages, but can be expressed by visibly pushdown languages [AM09]. These latter languages rely on a partitioning of the input alphabet into internal letters, call letters, and return letters, which naturally fits to program traces. A language over such a partitioned alphabet is *visibly pushdown* if there exists a *visibly pushdown automaton* (VPA) recognizing it. VPAs are pushdown automata, where operations on the stack are driven by the letter type: call letters can only push, return letters can only pop, while internal letters do not have access to the stack.

Runtime verification amounts to check a property during a program execution. In this paper we are interested in finding the earliest point during the execution where the property can be asserted, for properties defined by VPAs. More formally, given a word u over a partitioned alphabet and a VPA \mathcal{A} , we say that \mathcal{A} is *right-universal w.r.t. u* if, for every possible continuation v of u , the word uv is accepted by \mathcal{A} . Hence, if u is a prefix of the trace w and \mathcal{A} is known to be right-universal w.r.t. u , then it can be asserted that w verifies the property described by \mathcal{A} , without reading w entirely. Our aim is to check right-universality *incrementally* after each incoming event of the trace w , as described

Algorithm 1 Checking right-universality incrementally

```
function INCREMENTAL-RIGHT-UNIVERSALITY( $\mathcal{A}$ )  
   $u \leftarrow \epsilon$   
  while trace  $w$  is not completely read do  
     $a \leftarrow$  next letter of  $w$   
     $u \leftarrow ua$   
    if  $\mathcal{A}$  is right-universal w.r.t.  $u$  then  
      return True  
    end if  
  end while  
  return False  
end function
```

in Algorithm 1. By incremental, we mean that some information is propagated from one event (reading a letter in w) to the next one (reading the next letter in w), avoiding repeated identical computations. Note that our algorithm will not store u , but enough information for asserting right-universality of \mathcal{A} w.r.t. u .

For *safety properties*, right-universality allows to know at the earliest time point, whether the execution is sure, and thus whether controls can be stopped. This is typically the case for properties looking for a pattern (potentially complicated) that must occur during the execution. For *properties to be avoided*, this permits to stop the program before it enters an unsafe configuration, hence avoiding potential attacks [BJLW08].

These questions started to be addressed in the context of XML through *earliest query answering* of XPath expressions [BYFJ05,GNT09]. Indeed, XML documents also entail a nesting structure, similarly to program traces, and properties (or queries) over these documents can be expressed using VPAs. An algorithm asserting at the earliest time point, whether an XML document is accepted by a deterministic VPA has been proposed in [GNT09]. This algorithm runs in polynomial time at each incoming event. For non-deterministic VPAs, however, the decision problem associated with right-universality (i.e. given \mathcal{A} and u , is \mathcal{A} right-universal w.r.t. u ?) is known to be ExpTime-complete [GNT09].

Algorithms that incrementally check right-universality of *non-deterministic* VPAs are challenging in several aspects. First, VPAs are usually an intermediate object, resulting from the translation of a property expressed in some logics, like XPath for XML documents. Such translations rely on non-determinism, as for instance when referring to any call inside a given procedure call (which corresponds to the descendant axis in XPath). Any VPA can be determinized, but the procedure yields VPAs of exponential size [AM09]. Second, right-universality w.r.t. u can be considered as a variant of universality, parameterized by u . Recent techniques have been proposed to check universality of non-deterministic VPAs efficiently [TO12,FKL13,BDG13]. Among them, *antichains* have been successfully applied to decision problems related to non-deterministic automata: universality and inclusion for finite word automata [DDHR06], and for non-deterministic bottom-up tree automata [BHH⁺08]. Whether these techniques also apply to in-

crementally check right-universality is an interesting issue. Third, the algorithm proposed in the deterministic case [GNT09] does not directly generalize to the non-deterministic case.

Our contributions are the following. We propose an efficient incremental algorithm checking right-universality of a non-deterministic VPA \mathcal{A} . This algorithm relies on the progressive computation of *safe sets of configurations*, the configuration of a VPA being a state together with a stack content. It avoids to entirely determinize the given VPA. This algorithm uses antichains in order to get a compact representation of safe sets of configurations. We also propose efficient operators and data structures for saving useless computation. We report on some experiments, performed on randomly generated VPAs, and also on VPAs resulting from a translation from XPath expressions. These results exhibit the benefits of our optimizations.

Verification of traces with a nesting structure has already been addressed, through different aspects. Let us mention some of them. In [AEM04], the logic CaRet over words with a nesting structure is introduced, and a model checking algorithm is proposed. An extension is presented in [RCB08], with a corresponding monitor synthesis algorithm. VPAs are also sometimes used as an intermediate model to express properties on program traces [CA07,FJJ⁺12].

The paper is structured as follows. In Section 2 we introduce VPAs and right-universality. Safe sets of configurations are defined in Section 3. Our antichain-based algorithm is described in Section 4, and further optimizations in Section 5. Experiments are reported in Section 6. Proofs are provided in Appendix.

2 Visibly Pushdown Automata and Right-Universality

Visibly pushdown automata (VPAs) are pushdown automata working on a partitioned alphabet where only call symbols can push, return symbols can pop, and internal symbols can fire transitions without considering the stack [AM04,AM09].

In this paper, we consider trees, instead of words with a nesting structure (i.e. their linearization). Such a mapping is illustrated in Figure 1. Each call and its matching return are mapped to a node, and the calls and returns directly nested under this call correspond to the children of this node. These trees are unranked, as the number of children of a node is not determined by its label. We use VPAs as *unranked tree acceptors*, operating on their linearization [GNR08]. In particular these VPAs do not use internal symbols.

2.1 Unranked Trees

We here recall the standard definition of unranked trees, as provided for instance in [CDG⁺07]. Let Σ be a finite *alphabet*, and Σ^* (resp. Σ^+) be the set of all words (resp. non empty words) over Σ . The empty word is denoted by ϵ . Given two words $v, w \in \Sigma^*$ over Σ , v is a *prefix* (resp. *proper prefix*) of w if there exists a word $v' \in \Sigma^*$ (resp. $v' \in \Sigma^+$) such that $vv' = w$.

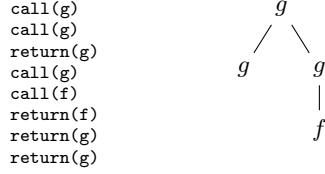


Fig. 1: Representation of a program trace as a tree t .

An *unranked tree* t over Σ is a tree such that its nodes are labeled by a letter of Σ and have an arbitrary number of children (the children are ordered from left to right). We call *a-node* a node with label $a \in \Sigma$. The set of all unranked trees over Σ is denoted by T_Σ . A *hedge* h over Σ is a finite sequence (empty or not) of unranked trees over Σ . The empty hedge is denoted by ϵ , and the set of all hedges over Σ is denoted by H_Σ .

Trees can be described by well-balanced words which correspond to a depth-first traversal of the tree. An opening tag is used to notice the arrival on a node and a closing tag to notice the departure from a node. For each $a \in \Sigma$, let a itself represent the opening tag and \bar{a} the related closing tag. The *linearization* $[t]$ of $t \in T_\Sigma$ is the *well-balanced* word over $\Sigma \cup \bar{\Sigma}$, with $\bar{\Sigma} = \{\bar{a} \mid a \in \Sigma\}$, inductively defined by: $[t] = a [t_1] \cdots [t_n] \bar{a}$, with a is the label of the root and t_1, \dots, t_n are its n subtrees (from left to right). The linearization is extended to hedges as follows. Let $h = t_1 \cdots t_n$ be the sequence of trees t_i , $1 \leq i \leq n$. Then $[h] = [t_1] \cdots [t_n]$. We denote by $[T_\Sigma]$ (resp. $[H_\Sigma]$) the set of linearizations of all trees in T_Σ (resp. hedges in H_Σ). Consider for instance the tree t in Figure 1: its linearization is the word $[t] = gg\bar{g}gf\bar{f}g\bar{g}$. Let $\text{Pref}(T_\Sigma)$ denote the set of all prefixes of $[T_\Sigma]$: $\text{Pref}(T_\Sigma) = \{u \in (\Sigma \cup \bar{\Sigma})^* \mid \exists v \in (\Sigma \cup \bar{\Sigma})^*, uv \in [T_\Sigma]\}$.

2.2 Visibly Pushdown Automata

Definition 1. A visibly pushdown automaton \mathcal{A} over a finite alphabet Σ is a tuple $\mathcal{A} = (Q, \Sigma, \Gamma, Q_i, Q_f, \Delta)$ where Q is a finite set of states containing initial states $Q_i \subseteq Q$ and final states $Q_f \subseteq Q$, a finite set Γ of stack symbols, and a finite set Δ of rules. Each rule in Δ is of the form $q \xrightarrow{a:\gamma} q'$ with $a \in \Sigma \cup \bar{\Sigma}$, $q, q' \in Q$, and $\gamma \in \Gamma$.

The left-hand side of a rule $q \xrightarrow{a:\gamma} p \in \Delta$ is (q, a) if $a \in \Sigma$, and (q, a, γ) if $a \in \bar{\Sigma}$. A VPA is *deterministic* if it has at most one initial state, and it does not have two distinct rules with the same left-hand side. We provide an example of deterministic VPA in Figure 2a.

A *configuration* of a VPA \mathcal{A} is a pair (q, σ) where $q \in Q$ is a state and $\sigma \in \Gamma^*$ a stack content. A configuration is *initial* (resp. *final*) if $q \in Q_i$ (resp. $q \in Q_f$) and $\sigma = \epsilon$. For $a \in \Sigma \cup \bar{\Sigma}$, we write $(q, \sigma) \xrightarrow{a} (q', \sigma')$ if there is a transition $q \xrightarrow{a:\gamma} q'$ in Δ verifying $\sigma' = \gamma \cdot \sigma$ if $a \in \Sigma$, and $\sigma = \gamma \cdot \sigma'$ if $a \in \bar{\Sigma}$. We extend this notation to words $u = a_1 \cdots a_n$, by writing $(q_0, \sigma_0) \xrightarrow{u} (q_n, \sigma_n)$

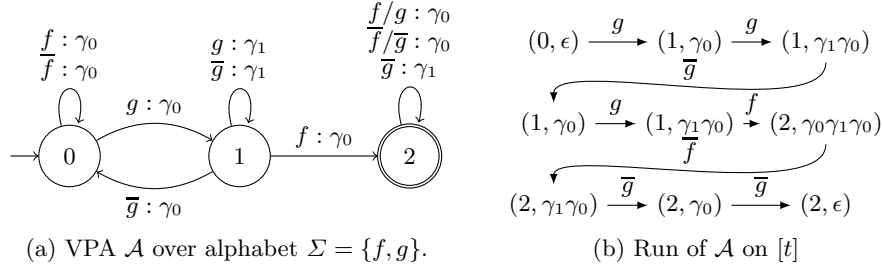


Fig. 2: An automaton and one of its runs.

whenever there exist configurations (q_i, σ_i) such that $(q_{i-1}, \sigma_{i-1}) \xrightarrow{a_i} (q_i, \sigma_i)$ for all $1 \leq i \leq n$.

A *run* of a VPA \mathcal{A} on a linearization $[t] = a_1 \cdots a_n$ of $t \in T_\Sigma$ is a sequence $(q_0, \sigma_0) \cdots (q_n, \sigma_n)$ of configurations such that (q_0, σ_0) is initial, and for every $1 \leq i \leq n$, $(q_{i-1}, \sigma_{i-1}) \xrightarrow{a_i} (q_i, \sigma_i)$. Such a run is *accepting* if (q_n, σ_n) is final. A tree t is *accepted* by \mathcal{A} if there is an accepting run on its linearization $[t]$. The set of accepted trees is called the *language* of \mathcal{A} and is written $L(\mathcal{A})$.³ For instance, given a tree t , the VPA of Figure 2a checks whether t has a g -node with an f -child. An accepting run on $[t]$ for the tree t of Figure 1, is depicted in Figure 2b.

2.3 Universality and Right-Universality

We conclude the section of preliminaries with the notions of universality⁴ and right-universality, that we will study in the remainder of the paper.

Definition 2. A VPA \mathcal{A} over Σ is said *universal* if \mathcal{A} accepts all trees $t \in T_\Sigma$, i.e. $L(\mathcal{A}) = T_\Sigma$. Let $u \in \text{Pref}(T_\Sigma) \setminus \{\epsilon\}$ be a non empty prefix of $[t_0]$ for some tree $t_0 \in T_\Sigma$. The VPA \mathcal{A} is said *right-universal w.r.t. u* if for all trees $t \in T_\Sigma$, if u is a prefix of $[t]$, then t is accepted by \mathcal{A} .

In other words, right-universality w.r.t. u allows to assert that any tree linearization beginning with u is accepted by the automaton. In this definition, notice that when $u = [t_0]$, then \mathcal{A} is right-universal w.r.t. u iff t_0 is accepted by \mathcal{A} . Moreover, as a universal VPA is right-universal w.r.t. all non empty words $u \in \text{Pref}(T_\Sigma)$, we assume in the sequel that VPAs are not universal.

In this article, our objective is to propose an *incremental* algorithm for right-universality of a VPA \mathcal{A} , in the sense described in the introduction (see Algorithm 1). More precisely, the linearization $[t_0]$ of a given tree t_0 is read letter

³ VPAs considered in this article are tree acceptors, with acceptance on empty stack.

⁴ This notion of universality is different from the one proposed for usual VPAs, where a VPA is universal if it accepts all words of $(\Sigma \cup \overline{\Sigma})^*$ (and not only linearizations of trees).

by letter⁵, and while \mathcal{A} is not right-universal w.r.t. the current read prefix u of $[t_0]$, the next letter of $[t_0]$ is read. When processing a new letter, we try to reuse prior computations as much as possible. In the sequel t_0 always refers to this particular tree.

We recall that the right-universality problem is ExpTime-complete for VPAs, but in PTime for deterministic VPAs [GNT09]. A dual problem to right-universality w.r.t. u is to ask whether every word starting with u is *rejected* by a given VPA. This problem is in PTime, even when the VPA is not deterministic. Indeed, this amounts to check whether no configuration reached after reading u can reach a final configuration.

2.4 Towards an Algorithm

In this section we consider several approaches for addressing right-universality of VPAs, and justify our choice of considering safe sets of configurations. Readers not familiar with the related literature can skip this section at first reading. As explained in the introduction, we discard the naive approach consisting in determinizing the VPA and then applying the algorithm in [GNT09], in order to avoid state-space explosion.

The algorithm given in [GNT09] for deterministic VPAs is a progressive computation of *safe states*, with the property that the VPA is right-universal w.r.t. u iff the state reached after reading u is safe. Note that safe states cannot be determined statically: a state can be safe at some position of the word, but not at another one. A first idea for the non-deterministic case is a subset construction, using sets of safe states instead of safe states, but this direct adaptation is not correct. Hence, *safe sets of configurations* will be the basis of our algorithm, instead of safe states. This point raises new questions, as the configuration space is infinite, unlike the state space. In fact we will see that at each event, safe sets of configurations have the same stack height as the depth of the word u leading to this event, and are thus of finite (but potentially huge) cardinality at each time point.

Recently, several authors have proposed efficient algorithms to check universality of VPAs [TO12,FKL13,BDG13]. The method that we give in [BDG13] computes in an incremental way the set \mathcal{R} of accessibility relations for all hedges (i.e. $\{\text{rel}_h \mid h \in H_\Sigma\}$ in the present paper, also called summaries in [AM04]). It uses *antichains* to get smaller objects to manipulate (in particular to limit \mathcal{R} to a strict smaller subset) and to avoid an explicit determinization step. It can be easily adapted for checking right-universality of VPAs [BDG12, Section 3.4]. However we face the problem of computing the whole set \mathcal{R} (instead of a strict subset): experiments indicate that this step is too much time consuming [BDG13]. For this reason we do not follow this approach, but use antichains over safe sets of configurations rather than accessibility relations. This approach is developed in the next section.

⁵ $[t_0]$ is the trace w of Algorithm 1.

3 Safe Sets of Configurations

In this section, $[t_0]$ is a fixed tree whose linearization is read letter by letter, and u denotes the current read prefix. We introduce the new notion of safe set of configurations related to u . We show how safe sets of configurations can be defined from such sets associated to smaller prefixes u' of $[t_0]$. This property will be useful to derive an incremental algorithm for checking right-universality of a VPA.

3.1 A Notion of Safety for Right-Universality

Definition 3. Let \mathcal{A} be a VPA and $\mathcal{C} \subseteq Q \times \Gamma^*$ be a non empty set of configurations. Let $u \in \text{Pref}(T_\Sigma) \setminus \{\epsilon\}$ be a non empty prefix of t_0 .

- \mathcal{C} is safe for u if for every v such that $uv \in [T_\Sigma]$, there exist $(q, \sigma) \in \mathcal{C}$ and $p \in Q_f$ such that $(q, \sigma) \xrightarrow{v} (p, \epsilon)$ in \mathcal{A} .
- \mathcal{C} is leaf-safe for u if for every $v = \bar{a}v'$ with $\bar{a} \in \bar{\Sigma}$ such that $uv \in [T_\Sigma]$, there exist $(q, \sigma) \in \mathcal{C}$ and $p \in Q_f$ such that $(q, \sigma) \xrightarrow{v} (p, \epsilon)$ in \mathcal{A} .⁶

Let $\text{Safe}(u) = \{\mathcal{C} \mid \mathcal{C} \text{ is safe for } u\}$ and $\text{LSafe}(u) = \{\mathcal{C} \mid \mathcal{C} \text{ is leaf-safe for } u\}$.

Intuitively, as stated in Theorem 1 below, if \mathcal{C} is the set of configurations reached in \mathcal{A} after reading u , then \mathcal{A} is right-universal w.r.t. u iff \mathcal{C} is safe for u . Indeed, for every possible v , one can find in \mathcal{C} at least one configuration leading to an accepting configuration after reading v . Before stating this theorem, let us give the next definition. Let $\text{Reach}(u)$ denote the set of configurations (q, σ) such that $(q_0, \sigma_0) \xrightarrow{u} (q, \sigma)$ for some initial configuration (q_0, σ_0) of \mathcal{A} .

Theorem 1. \mathcal{A} is right-universal w.r.t. u iff $\text{Reach}(u) \in \text{Safe}(u)$.

Let us illustrate this on the VPA \mathcal{A} depicted in Figure 2a. Recall that this VPA checks that the tree has a g -node with an f -child. After opening a g -root, the set of safe configurations is $\text{Safe}(g) = \{(q_2, \gamma_1), (q_2, \gamma_0)\}$, which does not contain the set of reached configurations, as $\text{Reach}(g) = \{(q_1, \gamma_0)\}$. Indeed, \mathcal{A} is not right-universal w.r.t. g , as no f -node has been encountered yet under a g -node. For the word gf , as expected, the situation differs: $\text{Safe}(gf) = \{(q_2, \gamma_1\gamma_0), (q_2, \gamma_0\gamma_0)\}$, and $\text{Reach}(gf) = \{(q_2, \gamma_0\gamma_0)\}$ is a safe set of configurations. By Theorem 1, \mathcal{A} is right-universal w.r.t. gf .

Hence, an algorithm computing both $\text{Reach}(u)$ and $\text{Safe}(u)$ can decide right-universality w.r.t. u . The set $\text{Reach}(u)$ is easy to compute, just by firing transitions of the VPA. In Sections 3.2-3.4, we detail how the set $\text{Safe}(u)$ can be defined from the set $\text{Safe}(u')$ with u' a proper prefix of u . In this way, while reading the linearization $[t_0]$ of a given tree t_0 , the set $\text{Safe}(u)$ with $u \neq \epsilon$ prefix of $[t_0]$, can be incrementally defined. In Section 4, we turn this approach into an algorithm.

⁶ The name “leaf-safe” comes from the fact that we only consider suffixes starting with a $\bar{\Sigma}$ symbol, and thus continuations of the tree that do not add children to the current node.

3.2 Starting Point

We begin with $\text{Safe}(a)$ with $a \in \Sigma$, for which we recall the definition.

$$\text{Safe}(a) = \{\mathcal{C} \subseteq Q \times \Gamma^* \mid \forall h \in H_\Sigma, \exists q_f \in Q_f, \exists (q, \sigma) \in \mathcal{C} : (q, \sigma) \xrightarrow{h\bar{a}} (q_f, \epsilon)\} \quad (1)$$

3.3 Reading a Letter $\bar{a} \in \bar{\Sigma}$

Suppose that we are reading an $\bar{a} \in \bar{\Sigma}$ and we have to compute $\text{Safe}(u\bar{a})$. Two cases occur : either $u\bar{a} \neq [t_0]$, or $u\bar{a} = [t_0]$. Let us study these two cases and show how to manage them from an algorithmic point of view.

If $u\bar{a} \neq [t_0]$, we can retrieve safe sets configurations from prior computed such sets: $\text{Safe}(u\bar{a}) = \text{Safe}(u')$ where $u' \neq \epsilon$ is the unique prefix of u such that $u = u'a[h]$, with $h \in H_\Sigma$. Indeed as shown by Lemma 1 below, we have $\text{Safe}(u'a[h]\bar{a}) = \text{Safe}(u')$. Hence, from an algorithmic point of view, we just have to use a stack \mathcal{S} to store these safe sets of configurations. When opening a , we put $\text{Safe}(u')$ on the stack, and when closing \bar{a} , we pop it. As h is a hedge, the stack before reading \bar{a} is exactly the stack after reading a .

Lemma 1. *If $h \in H_\Sigma$, then $\text{Safe}(u[h]) = \text{Safe}(u)$ and $L\text{Safe}(u[h]) = L\text{Safe}(u)$.*

If $u\bar{a} = [t_0]$, the previous argument is no longer correct since $u' = \epsilon$ and $\text{Safe}(u')$ is not defined for the empty word. Nevertheless, by definition $\text{Safe}(u\bar{a}) = \{\mathcal{C} \subseteq Q \times \Gamma^* \mid \exists (q, \epsilon) \in \mathcal{C} \text{ with } q \in Q_f\}$. Therefore we can again use stack \mathcal{S} as described before if it is initialized with the set

$$\text{Init} = \{\mathcal{C} \subseteq Q \times \Gamma^* \mid \exists (q, \epsilon) \in \mathcal{C} \text{ with } q \in Q_f\}. \quad (2)$$

3.4 Reading a Letter $a \in \Sigma$

Let us now consider the much more involved case of sets $\text{Safe}(ua)$ with $a \in \Sigma$. When reading an $a \in \Sigma$, two successive steps are performed, with leaf-safe sets of configurations as intermediate object:

$$\text{Safe}(u) \xrightarrow{\text{Step 1}} L\text{Safe}(ua) \xrightarrow{\text{Step 2}} \text{Safe}(ua)$$

For this purpose, we introduce the notion of *predecessor*. Let $\mathcal{C}, \mathcal{C}'$ be two sets of configurations, let $\bar{a} \in \bar{\Sigma}$ and $h \in H_\Sigma$.

- \mathcal{C} is an \bar{a} -predecessor of \mathcal{C}' if $\forall (q', \sigma') \in \mathcal{C}', \exists (q, \sigma) \in \mathcal{C}, (q, \sigma) \xrightarrow{\bar{a}} (q', \sigma')$.
- \mathcal{C} is a h -predecessor of \mathcal{C}' if $\forall (q', \sigma') \in \mathcal{C}', \exists (q, \sigma) \in \mathcal{C}, (q, \sigma) \xrightarrow{[h]} (q', \sigma')$.

Let $\text{Pred}_{\bar{a}}(\mathcal{C}') = \{\mathcal{C} \mid \mathcal{C} \text{ is an } \bar{a}\text{-predecessor of } \mathcal{C}'\}$ and $\text{Pred}_h(\mathcal{C}') = \{\mathcal{C} \mid \mathcal{C} \text{ is a } h\text{-predecessor of } \mathcal{C}'\}$.

The next proposition can be used to perform Step 1 and Step 2. It states that safe sets of configurations are only among predecessors of prior safe sets of configurations.

Proposition 1. *Let $ua \in \text{Pref}(T_\Sigma)$.*

$$\mathcal{C} \in \text{LSafe}(ua) \iff \exists \mathcal{C}' \in \text{Safe}(u), \mathcal{C} \in \text{Pred}_{\bar{a}}(\mathcal{C}') \quad (3)$$

$$\mathcal{C} \in \text{Safe}(ua) \iff \forall h \in H_\Sigma, \exists \mathcal{C}' \in \text{LSafe}(ua), \mathcal{C} \in \text{Pred}_h(\mathcal{C}') \quad (4)$$

However, to get an algorithm, we face the problem that the number of hedges to consider in Equivalence (4) is infinite. We use relations to overcome this. Also the size of $\text{Safe}(u)$ may be huge and not all configurations of $\text{Safe}(u)$ are crucial for checking right-universality w.r.t. u . We use antichains to get a compact representation of $\text{Safe}(u)$. These two concepts are explained in the following section.

4 An Antichain-Based Algorithm for Right-Universality

In this section, we give an antichain-based algorithm for incrementally checking right-universality of a VPA, that uses the approach given in Sections 3.2-3.4.

Let us summarize this approach. Let $[t_0]$ be the linearization of a given tree t_0 . We initialize a stack \mathcal{S} with set Init defined in (2) and we start by computing $\text{Safe}(u)$ with u being the first letter of $[t_0]$ (see (1)). Suppose that $\text{Safe}(u)$ has been computed from the current read prefix u of $[t_0]$. Then, if the next letter read in $[t_0]$ is $a \in \Sigma$, we compute $\text{Safe}(ua)$ from $\text{Safe}(u)$ by Steps 1 and 2, and we put $\text{Safe}(ua)$ on the stack \mathcal{S} . If the next letter read in $[t_0]$ is $\bar{a} \in \bar{\Sigma}$, then we pop the stack \mathcal{S} . The element that has been popped is the set $\text{Safe}(u') = \text{Safe}(u\bar{a})$ where u' is the unique prefix of u such that $u = u'a[h]$ (except if $u\bar{a} = t_0$ in which case the popped set is equal to $\text{Init} = \text{Safe}(u\bar{a})$), see Section 3.3. At each computation of $\text{Safe}(u)$, we check whether $\text{Reach}(u) \in \text{Safe}(u)$ (see Theorem 1).

4.1 Finite Number of Hedges

We begin by showing that only a finite number of hedges has to be considered in Equivalence (4). The reason is that a hedge h does not change the stack of a configuration during a run of a VPA, that is $(q, \sigma) \xrightarrow{[h]} (q', \sigma)$. So h can be considered as a function mapping each state q to the set of states obtained when traversing h from q . Formally, for every $h \in H_\Sigma$, rel_h is the function from Q to 2^Q such that $q' \in \text{rel}_h(q)$ iff $(q, \sigma) \xrightarrow{[h]} (q', \sigma)$ for some $\sigma \in \Gamma^*$. The number of such functions is finite, and bounded by $|Q| \cdot 2^{|Q|}$. These functions naturally define an equivalence relation of finite index over H_Σ :

$$h \sim h' \iff \text{rel}_h = \text{rel}_{h'}.$$

Let us note H for a subset containing one hedge per \sim -class. We have $|H| \leq |Q| \cdot 2^{|Q|}$. The next lemma indicates that the computation of h -predecessors can be limited to $h \in H$.

Lemma 2. *For every $h \in H_\Sigma$, \mathcal{C} is a h -predecessor of \mathcal{C}' iff there exists $h' \in H$ with $h \sim h'$, such that \mathcal{C} is a h' -predecessor of \mathcal{C}' .*

The set H can be computed by a saturation method based on the VPA rules, as described by Algorithm 3 in Appendix B.2.

4.2 Antichains

Let us now introduce antichains. Consider the set $2^{Q \times \Gamma^*}$ of all sets of configurations, with the \subseteq operator. An *antichain* is a set of pairwise incomparable sets of configurations with respect to \subseteq . Given a set α of sets of configurations, we denote by $\lfloor \alpha \rfloor$ the \subseteq -minimal elements of α . The set α is \subseteq -upward closed if for all $\mathcal{C} \in \alpha$ and $\mathcal{C}' \subseteq \mathcal{C}$, we have $\mathcal{C}' \in \alpha$. From their definition, it is immediate that $\text{Safe}(u)$, $\text{LSafe}(u)$, $\text{Pred}_{\bar{a}}(\mathcal{C})$ and $\text{Pred}_h(\mathcal{C})$ are \subseteq -upward closed sets.

The idea is to compute the antichain $\lfloor \text{Safe}(u) \rfloor$ (resp. $\lfloor \text{LSafe}(u) \rfloor$) instead of the whole \subseteq -upward closed set $\text{Safe}(u)$ (resp. $\text{LSafe}(u)$). The next corollary of Theorem 1 indicates that such a limited computation is enough to check whether a VPA \mathcal{A} is right-universal w.r.t. u .

Corollary 1. *\mathcal{A} is right-universal w.r.t. u iff there exists $\mathcal{C} \in \lfloor \text{Safe}(u) \rfloor$ such that $\mathcal{C} \subseteq \text{Reach}(u)$.*

Moreover, we show in Appendix B.3 that the antichains $\lfloor \text{Safe}(u) \rfloor$ and $\lfloor \text{LSafe}(u) \rfloor$ are finite and only contain finite sets \mathcal{C} of configurations such that $\mathcal{C} \subseteq Q \times \Gamma^k$ for some k . We now try to use these antichains at the starting point, and in Steps 1 and 2 of our approach.

4.3 Starting Point with Antichains

Let us explain how to compute $\text{Safe}(a)$. Clearly, by definition of H and using (1) (see Section 3.2), we can compute $\lfloor \text{Safe}(a) \rfloor$ as follows:

$$\lfloor \text{Safe}(a) \rfloor = \left\lfloor \left\{ \mathcal{C} \mid \forall h \in H, \exists q_f \in Q_f, \exists (q, \sigma) \in \mathcal{C} : (q, \sigma) \xrightarrow{h\bar{a}} (q_f, \epsilon) \right\} \right\rfloor. \quad (5)$$

4.4 Step 1 with Antichains: from $\lfloor \text{Safe}(u) \rfloor$ to $\lfloor \text{LSafe}(ua) \rfloor$

For the two steps, the goal is to adapt Proposition 1 so that it uses $\lfloor \text{Safe}(\cdot) \rfloor$ instead of $\text{Safe}(\cdot)$, and $\lfloor \text{LSafe}(\cdot) \rfloor$ instead of $\text{LSafe}(\cdot)$. We begin with Step 1. Implication (\Rightarrow) of Equivalence (3) can be directly adapted.

$$\mathcal{C} \in \lfloor \text{LSafe}(ua) \rfloor \implies \exists \mathcal{C}' \in \lfloor \text{Safe}(u) \rfloor, \mathcal{C} \in \text{Pred}_{\bar{a}}(\mathcal{C}'). \quad (6)$$

Implication (6) gives us a way to compute $\lfloor \text{LSafe}(ua) \rfloor$ from $\lfloor \text{Safe}(u) \rfloor$: it suffices to take all \bar{a} -predecessors of elements of $\lfloor \text{Safe}(u) \rfloor$ and then limit to those predecessors that are \subseteq -minimal. We can even only consider minimal \bar{a} -predecessors of $\lfloor \text{Safe}(u) \rfloor$ in the following sense: \mathcal{C} is a *minimal \bar{a} -predecessor* of \mathcal{C}' if for all \mathcal{C}'' \bar{a} -predecessor of \mathcal{C}' , $\mathcal{C}'' \subseteq \mathcal{C} \implies \mathcal{C}'' = \mathcal{C}$. We finally obtain:

$$\lfloor \text{LSafe}(ua) \rfloor = \left\{ \mathcal{C} \mid \mathcal{C} \text{ is a minimal } \bar{a}\text{-predecessor of } \mathcal{C}' \in \lfloor \text{Safe}(u) \rfloor \right\}. \quad (7)$$

4.5 Step 2 with Antichains: from $\llbracket L\text{Safe}(ua) \rrbracket$ to $\llbracket \text{Safe}(ua) \rrbracket$

The second step for computing $\llbracket \text{Safe}(ua) \rrbracket$ from $\llbracket \text{Safe}(u) \rrbracket$ relies on the introduction of antichains in Equivalence (4). Implication (\Rightarrow) holds with antichains:

$$\mathcal{C} \in \llbracket \text{Safe}(ua) \rrbracket \implies \forall h \in H_\Sigma, \exists \mathcal{C}' \in \llbracket L\text{Safe}(ua) \rrbracket, \mathcal{C} \in \text{Pred}_h^-(\mathcal{C}'). \quad (8)$$

Similarly to Implication (6), we can restrict h -predecessors to only consider minimal ones: \mathcal{C} is a *minimal h -predecessor* of \mathcal{C}' if for all \mathcal{C}'' h -predecessor of \mathcal{C}' , $\mathcal{C}'' \subseteq \mathcal{C} \implies \mathcal{C}'' = \mathcal{C}$. Moreover, by Lemma 2, we can limit the computations to hedges $h \in H$ where H is the finite set introduced in Section 4.1. Therefore, we obtain the next equality.

$$\llbracket \text{Safe}(ua) \rrbracket = \left[\left\{ \mathcal{C} \mid \mathcal{C} = \bigcup_{h \in H} \mathcal{C}_h \text{ with } \mathcal{C}_h \text{ a minimal } h\text{-predecessor of } \mathcal{C}' \in \llbracket L\text{Safe}(ua) \rrbracket \right\} \right] \quad (9)$$

4.6 Algorithm

We are now able to give our antichain-based algorithm (see Algorithm 2). It uses a stack \mathcal{S} (initially empty) as recalled at the beginning of this section. The computation of $\llbracket H \rrbracket$ is considered as a *preprocessing*, as its value only depends on \mathcal{A} and not on $[t_0]$. The used results are mentioned inside the algorithm.

5 Improvements and Implementation

Section 4 resulted in a first algorithm for incrementally testing whether a VPA is right-universal. In this section, we show how this algorithm can be improved by limiting hedges to consider, and optimizing operators and predecessors to be computed. We also give the improved algorithm in a more detail way than in Algorithm 2, as well as the underlying data structures.

5.1 Minimal Hedges

A first improvement is obtained by further restricting hedges to consider. Indeed it suffices to consider *minimal hedges* wrt their function rel_h . Formally, let us write $h \leq h'$ whenever $\text{rel}_h(q) \subseteq \text{rel}_{h'}(q)$ for every $q \in Q$. We denote by $\llbracket H \rrbracket$ the \leq -minimal elements of H . From the definition of h -predecessor, we have:

$$\mathcal{C} \text{ } h\text{-predecessor of } \mathcal{C}' \text{ and } h \leq h' \implies \mathcal{C} \text{ } h'\text{-predecessor of } \mathcal{C}' \quad (10)$$

This property can be used to replace $h \in H$ in (9) by $h \in \llbracket H \rrbracket$:

$$\llbracket \text{Safe}(ua) \rrbracket = \left[\left\{ \mathcal{C} \mid \mathcal{C} = \bigcup_{h \in \llbracket H \rrbracket} \mathcal{C}_h \text{ with } \mathcal{C}_h \text{ a minimal } h\text{-predecessor of } \mathcal{C}' \in \llbracket L\text{Safe}(ua) \rrbracket \right\} \right] \quad (11)$$

Algorithm 2 Checking right-universality incrementally with antichains

```

function ANTICHAIN-BASED INCREMENTAL-RIGHT-UNIVERSALITY( $\mathcal{A}$ ,  $\lfloor H \rfloor$ )
  PUSH( $\mathcal{S}$ ,  $\lfloor Init \rfloor$ )  % by (2),  $\mathcal{S}$  being a stack
   $u \leftarrow$  first letter of  $\lfloor t_0 \rfloor$ 
   $R \leftarrow$   $\langle$ Compute  $Reach(u)$  $\rangle$ 
   $\alpha \leftarrow$   $\langle$ Compute  $\lfloor Safe(u) \rfloor$  $\rangle$   % by (5)
  if  $\exists \mathcal{C} \in \alpha : \mathcal{C} \subseteq R$  then
    return True  % Corollary 1
  end if
  while  $\lfloor t_0 \rfloor$  is not completely read do
     $a \leftarrow$  next letter of  $\lfloor t_0 \rfloor$ 
     $u \leftarrow ua$ 
     $R \leftarrow$   $\langle$ Compute  $Reach(u)$  from  $R$  $\rangle$ 
    if  $a \in \overline{\mathcal{S}}$  then
       $\alpha \leftarrow$  POP( $\mathcal{S}$ )
    else
      PUSH( $\mathcal{S}$ ,  $\alpha$ )
       $\alpha \leftarrow$   $\langle$ Compute  $\lfloor Safe(u) \rfloor$  from  $\alpha$  $\rangle$   % by (7) and (9)
    end if
    if  $\exists \mathcal{C} \in \alpha : \mathcal{C} \subseteq R$  then
      return True  % by Corollary 1
    end if
  end while
  return False
end function

```

and similarly for the starting point:

$$\lfloor Safe(a) \rfloor = \left[\left\{ \mathcal{C} \mid \forall h \in \lfloor H \rfloor, \exists q_f \in Q_f, \exists (q, \sigma) \in \mathcal{C} : (q, \sigma) \xrightarrow{h\bar{a}} (q_f, \epsilon) \right\} \right] \quad (12)$$

5.2 An Appropriate Operator

Equation (11) expresses that every set of configurations \mathcal{C} in $\lfloor Safe(ua) \rfloor$ is the union of \mathcal{C}_h with $h \in \lfloor H \rfloor$. We introduce a new operator to improve the readability and find new properties. Let S be a finite set, and $A, B \in 2^{2^S \setminus \{\emptyset\}}$. The set $A \sqcup B \in 2^{2^S}$ is defined by

$$A \sqcup B = \{a \cup b \mid a \in A \text{ and } b \in B\}.$$

Operator \sqcup builds sets obtained by taking one set of each of its operands, and performing their union. It is obviously associative and commutative. Notice that the elements of A, B are supposed to be non empty sets. This is always the case in the algorithms using this operator.

When combined with operator $\lfloor \cdot \rfloor$, operands of the \sqcup operator can be splitted, so that \sqcup is to be computed on smaller sets:

$$\lfloor A \sqcup B \rfloor = \lfloor (A \cap B) \cup (A \setminus B \sqcup B \setminus A) \rfloor \quad (13)$$

We experimentally observed that a good strategy to evaluate an expression $A_1 \sqcup \dots \sqcup A_n$ is to process sets A_i with increasing cardinality.

Equation (11) can now be rewritten in a simpler way as follows.

$$\lfloor \text{Safe}(ua) \rfloor = \left[\bigsqcup_{h \in [H]} \{ \mathcal{C}_h \mid \mathcal{C}_h \text{ is a minimal } h\text{-predecessor of } \mathcal{C}' \in [L\text{Safe}(ua)] \} \right] \quad (14)$$

We can go further by reconsidering the notions of minimal \bar{a} - and h -predecessor with the new operator \sqcup . From the definition of minimal predecessor, we immediately get that \mathcal{C} is a minimal \bar{a} -predecessor of \mathcal{C}' iff \mathcal{C} belongs to the set $\left[\bigsqcup_{(q', \sigma') \in \mathcal{C}'} \{ \{(q, \sigma)\} \mid (q, \sigma) \xrightarrow{\bar{a}} (q', \sigma') \} \right]$, and similarly for h -minimal predecessors. We provide in Appendix C.3 a similar improvement for $\lfloor \text{Safe}(a) \rfloor$.

5.3 Using a SAT Solver to Find Minimal Predecessors

The operator \sqcup allows to compute Step 1 and Step 2. Equation (13) accelerates these computations. In this section, we propose a method to compute Step 1 based on a SAT solver (a similar approach also works for Step 2, see Appendix C.4).

A SAT solver is an algorithm used to efficiently test the satisfiability of a boolean formula φ , that is, to check whether there exists a valuation v of the boolean variables of φ that makes φ true. In this case we say that v is a *model* of φ , denoted by $v \models \varphi$.

In Step 1, the computation of set $\lfloor L\text{Safe}(ua) \rfloor$ from $\lfloor \text{Safe}(u) \rfloor$ is given in (7):

$$\lfloor L\text{Safe}(ua) \rfloor = \{ \mathcal{C} \mid \mathcal{C} \text{ is an } \bar{a}\text{-predecessor of } \mathcal{C}' \in \lfloor \text{Safe}(u) \rfloor \}.$$

We recall that \mathcal{C} is an \bar{a} -predecessor of \mathcal{C}' if for all $(q', \sigma') \in \mathcal{C}'$, there exists $(q, \sigma) \in \mathcal{C}$ such that $(q, \sigma) \xrightarrow{\bar{a}} (q', \sigma')$. We also recall that $\mathcal{C} \in \lfloor L\text{Safe}(ua) \rfloor$ is a finite object such that $\mathcal{C} \subseteq Q \times \Gamma^k$ for some k . Let us associate a boolean variable x_c to each configuration $c \in Q \times \Gamma^k$.

We consider the following boolean formula $\varphi_{\bar{a}}$:

$$\varphi_{\bar{a}} = \bigvee_{\mathcal{C}' \in \lfloor \text{Safe}(u) \rfloor} \bigwedge_{c' \in \mathcal{C}'} \bigvee_{c \xrightarrow{\bar{a}} c'} x_c,$$

Let $v_{\mathcal{C}}$ be the valuation such that $v_{\mathcal{C}}(x_c) = 1$ iff $c \in \mathcal{C}$. We immediately obtain that:

$$v_{\mathcal{C}} \models \varphi_{\bar{a}} \quad \text{iff} \quad \mathcal{C} \in L\text{Safe}(ua) \cap Q \times \Gamma^k.$$

We define an ordering over valuations as follows, in a way to have a notion of minimal models equivalent to \subseteq -minimal elements of $L\text{Safe}(ua)$.

Let V be a set of boolean variables, let v and v' be two valuations over V . We define $v' \leq v$ iff for all variables $x \in V$, $v'(x) = 1 \implies v(x) = 1$. We denote

$v' < v$ if $v' \leq v$ and $v' \neq v$. Given φ a boolean formula over V , we say that a model v of φ is *minimal* if for all model v' of φ , we have $v' \leq v \implies v' = v$. We get the next characterization.

Lemma 3. $v_{\mathcal{C}}$ is a minimal model of $\varphi_{\bar{\alpha}}$ iff $\mathcal{C} \in \lfloor L\text{Safe}(ua) \rfloor$.

We can now explain how to compute all the minimal models of formula $\varphi_{\bar{\alpha}}$. Let φ be a boolean formula over V .

First, we explain, knowing a model v of φ , how to compute a model v' of φ such that $v' < v$ (if it exists). Consider the next formula φ' :

$$\varphi' = \varphi \wedge \left(\bigwedge_{x \in V_0} \neg x \right) \wedge \left(\bigvee_{x \in V_1} \neg x \right)$$

where V_0 (respectively V_1) is the set of all variables $x \in V$ such that $v(x) = 0$ (resp. $v(x) = 1$). If φ' has a model v' , it follows from the definition of φ' that v' is a model of φ such that $v' < v$. Otherwise, v is a minimal model of φ . So from a model of φ we can compute a minimal model of φ by repeating the above procedure.

Second, let us explain how to compute all the minimal models of φ . Suppose that we already know some minimal model v of φ , and let V_1 be the set of variables $x \in V$ such that $v(x) = 1$. Consider the formula

$$\varphi' = \varphi \wedge \left(\bigvee_{x \in V_1} \neg x \right).$$

Then a model v' of φ' , if it exists, is a model of φ such that neither $v' < v$ (since v is minimal) nor $v < v'$ (by definition of φ'). With the previous procedure, we thus get a minimal model of φ that is distinct from v . In this way we can compute all minimal models of φ .

5.4 Improved Algorithm and Data Structures

Let us come back to Algorithm 2 by indicating the underlying data structures and the improvements resulting from the previous three sections.

As explained in Section 5.1, we restrict the computations to the set $\lfloor H \rfloor$ of minimal hedges. Notice that Algorithm 3 that computes the set $\{\text{rel}_h \mid h \in H\}$ (see Appendix) can be easily adapted to compute the set of its \leq -minimal elements. In the main loop of Algorithm 2, Steps 1 and 2 can be computed either with the new operator \sqcup or with a SAT solver (see Sections 5.2 and 5.3 resp.).

Efficient data structures are used both for the relations associated to minimal hedges and for the antichains $\lfloor H \rfloor$, $\lfloor \text{Safe}(u) \rfloor$ and $\lfloor L\text{Safe}(u) \rfloor$. A relation rel_h , with $h \in \lfloor H \rfloor$, is stored as an array of bit-vectors. In this way the composition is computed efficiently using bit-operations, as well as the number of elements of the relation. A hash table is used to store each antichain, such that elements with different weights are stored in different lists. In the case of $\lfloor H \rfloor$, the weight is the number of elements of $\text{rel}_h \in \lfloor H \rfloor$. In the case of $\lfloor \text{Safe}(u) \rfloor$ (resp. $\lfloor L\text{Safe}(u) \rfloor$), the weight is the number of elements of $\mathcal{C} \in \lfloor \text{Safe}(u) \rfloor$ (resp. $\mathcal{C} \in \lfloor L\text{Safe}(u) \rfloor$). In this way, comparing a new element with the elements of the antichain is made more efficient, by limiting the comparison with elements of the same weight.

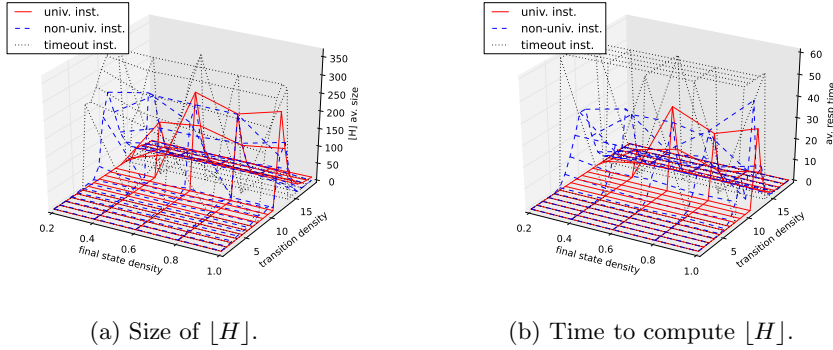


Fig. 3: Space and time consumption for computing $[H]$ on 50 VPAs of various densities, $|Q| = 10$, $|\Sigma| = |\Gamma| = 3$, and timeout of 60s.

6 Experiments

We have implemented Algorithm 2 in a prototype tool together with the data structures and improvements proposed in Section 5 (following both approaches, using operator \sqcup and a SAT solver). We mainly use Python, with C binding to the Glucose SAT solver (first ranked in recent SAT competitions) [Glucose].

The experimental tests are performed on two different benchmarks, one composed of randomly generated VPAs, and another one based on the translation of XPath expressions to VPAs. The experiments were run on a PC equipped with an Intel i7 2.8GHz processor, 6 GB of RAM and running Linux Ubuntu 3.2.

6.1 Randomly Generated VPAs

During the random generation of VPAs $\mathcal{A} = (Q, \Sigma, \Gamma, Q_i, Q_f, \Delta)$, $|Q_i|$ is fixed to 1, and several parameters vary: the sizes $|Q|$, $|\Sigma|$ and $|\Gamma|$, the density of final states $\frac{|Q_f|}{|Q|}$, and the transition density⁷. Our online algorithm for checking right-universality needs to compute the set $[H]$ of minimal hedges as a preprocessing. This set can be huge and thus lead to a timeout. In Figure 3a (resp. Figure 3b), we indicate the average size of $[H]$ (resp. the average execution time to compute it) for randomly generated VPAs with variable transition density (from 1 to 19) and variable final state density. In this test, we distinguish universal automata from non-universal ones since a universal VPA is right-universal w.r.t all $u \in Pref(T_\Sigma) \setminus \{\epsilon\}$. The two figures show that timeout happens for instances with transition density around 12, and that the density of final states has few influence.

⁷ equal to the number of outgoing transitions per state and per symbol.

$ Q $	universality		right-universality		null rel_h		timeout (preproc.)		timeout (online)	
	tr. density: 8	tr. density: 16	8	16	8	16	8	16	8	16
10	12 (0.59/0.01)	40 (1.19/0.01)	22 (16.62/9.36)	44 (11.00/1.10)	52 (0.67/0.01)	0	2	2	4	
20	13 (9.35/0.01)	29 (0.86/0.01)	14 (20.44/30.85)	41 (5.18/1.07)	46 (9.62/0.01)	0	14	17	3	3
30	9 (20.22/0.01)	34 (13.39/0.01)	11 (21.23/5.85)	34 (27.31/3.43)	40 (1.23/0.01)	0	30	21	0	1

Table 1: Number for each type (universal, right-universal, null rel_h , timeout) among 90 instances, with a timeout of 300s, and various number of states and transition densities. Average (preprocessing/online) time is given in parentheses (in s).

VPA id	Response	Time (s.)	$ H $	$Reach(u)$	$Safe(u)$
q20-a02-x02-o16-c16-f0.5-00	right-universal w.r.t. a_1	0.03	3	4	33
q20-a02-x04-o16-c16-f0.5-02	right-universal w.r.t. a_1	0.07	3	8	53
q20-a02-x04-o16-c16-f0.5-06	right-universal w.r.t. a_0	0.06	3	5	54
q20-a03-x02-o16-c16-f0.5-06	right-universal w.r.t. a_2	0.13	12	4	32
q20-a03-x03-o16-c16-f0.5-06	right-universal w.r.t. a_0	0.37	23	4	59
q20-a03-x03-o16-c16-f0.5-09	right-universal w.r.t. a_2a_0	1.74	13	14	179
q20-a03-x04-o16-c16-f0.5-07	right-universal w.r.t. a_2	1.02	44	5	116
q20-a04-x02-o16-c16-f0.5-03	right-universal w.r.t. a_3	0.71	71	5	75
q20-a04-x02-o16-c16-f0.5-07	right-universal w.r.t. a_0	1.07	81	5	74
q20-a04-x03-o16-c16-f0.5-03	right-universal w.r.t. a_2	15.88	359	5	447

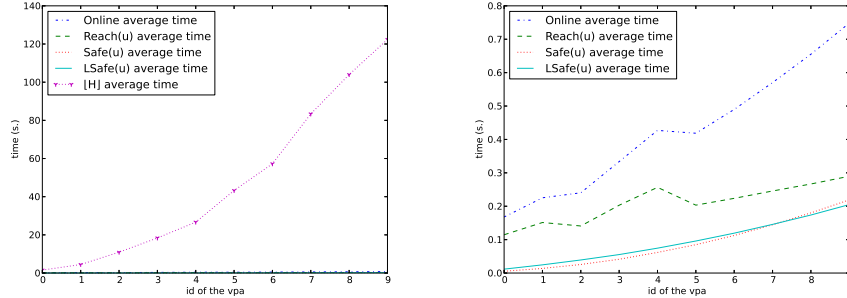
Table 2: Time and data structures size on random VPAs, with transition density of 16, $|Q| = 20$, $|\Sigma| = |\Gamma| \in \{2, 3, 4\}$, and timeout of 300s.

In the next experiment, we study the behavior of our algorithm on 90 random instances of size 10 (resp. 20, 30), with fixed transition density 8 (resp. 16) and fixed density 0.5 of final states⁸. The considered tree t_0 is a complete binary tree up to height 3 filled with randomly generated letters of Σ . We only comment the experiment using a SAT solver since it outperforms the approach with operator \sqcup . We observe several behaviors (see Table 1): many automata are either universal, or right-universal w.r.t. u with $|u| = 1$ (except 7 cases where $|u| = 2$), or exhibits a hedge h with rel_h being the null relation⁹; the number of timeout increases with $|Q|$. Table 2 indicates, for some representative VPAs, the execution time and the sizes of $|H|$, $Reach(u)$ and $|Safe(u)|$.

When it declares that a given VPA is right-universal w.r.t. u , our algorithm has the nice property that it reproduces the *same execution* (thus with the same time) for each tree t_0 such that u is prefix of $[t_0]$. This is clearly not the case for the *membership algorithm* that computes $Reach([t_0])$ and checks if it contains a final configuration. For instance, on a random VPA with size $|Q| = 20$ and transition density 8, our algorithm consumes less than 5s to declare that the

⁸ to deal with manageable and not too small sets $[H]$.

⁹ This means that the automaton is never right-universal w.r.t. u , for any proper prefix u of $[t_0]$. Therefore, in such cases, our algorithm is slower than the membership algorithm.



(a) With the $[H]$ average time curve. (b) Without the $[H]$ average time curve.

Fig. 4: Average time for computing $Safe(u)$, $LSafe(u)$ and $Reach(u)$ and overall computation for 15 random trees.

automaton is right-universal w.r.t. a , whereas the membership algorithm takes more than 300s as soon as the height of a binary tree t_0 with an a -root is equal to 8 (see also Appendix D.1).

6.2 VPAs resulting from XPath translation

Our second benchmark is based on VPAs obtained from queries over XML documents expressed in the XPath language, and then translated into VPAs. This translation was performed by the QuiXProc tool, as described in [GN11]. This family of XPath expressions yields VPAs of linear size increase (VPA with id i has $16 + 11i$ states), and looks for some complex patterns in the tree. In Figure 4 we report the time used by our algorithm on randomly generated trees, for this family of VPAs. This shows that for real-world VPAs, the size of $[H]$ is outside the hardness threshold exhibited in Figure 3. For instance, for the VPA with id 9 and size $|Q| = 115$, $[H]$ is computed in about 120s. Moreover, Figure 4b shows the efficiency of our online algorithm (less than 0.8s, see also Appendix D.2).

Acknowledgements

We thank Grégoire Sutre for useful discussions, and the QuiXProc team for translating XPath expressions of our benchmark to VPAs. We also thank one of the reviewers who helped us to improve the presentation of this article. The second author is supported by a grant from FRIA. This work was also partially supported by CNRS SOSP project.

References

- [AEM04] R. Alur, K. Etessami, and P. Madhusudan, *A temporal logic of nested calls and returns*, Proc. TACAS, LNCS **2988**, Springer, 2004, 67–481.

- [AM04] R. Alur and P. Madhusudan, *Visibly pushdown languages*, Proc. STOC, ACM-Press, 2004, pp. 202–211.
- [AM09] ———, *Adding nesting structure to words*, J. ACM **56** (2009), 1–43.
- [BYFJ05] Z. Bar-Yossef, M. Fontoura, and V. Josifovski, *Buffering in query evaluation over XML streams*, Proc. PODS, ACM-Press, 2005, 216–227.
- [BJLW08] M. Benedikt, A. Jeffrey, and R. Ley-Wild, *Stream Firewalling of XML Constraints*, Proc. SIGMOD Conference, ACM-Press, 2008, 487–498.
- [BHH⁺08] A. Bouajjani, P. Habermehl, L. Holk, T. Touili, and T. Vojnar, *Antichain-based universality and inclusion testing over nondeterministic finite tree automata*, Proc. CIAA, LNCS **5148**, Springer, 2008, 57–67.
- [BDG12] V. Bruyère, M. Ducobu, and O. Gauwin, *Visibly pushdown automata on trees: universality and u-universality*, CoRR abs/1205.2841, 2012.
- [BDG13] ———, *Visibly pushdown automata: Universality and inclusion via antichains*, Proc. LATA, LNCS **7810**, Springer, 2013, 190–201.
- [BYBC10] T. Bultan, F. Yu, and A. Betin-Can, *Modular verification of synchronization with reentrant locks*, Proc. MEMOCODE, IEEE Computer Society, 2010, 59–68.
- [CA07] S. Chaudhuri and R. Alur, *Instrumenting C programs with nested word monitors*, Proc. SPIN, LNCS **4595**, Springer, 2007, 279–283.
- [CDG⁺07] H. Comon, M. Dauchet, R. Gilleron, C. Löding, F. Jacquemard, D. Lugiez, S. Tison, and M. Tommasi, *Tree automata techniques and applications*, Available on: <http://www.grappa.univ-lille3.fr/tata>, 2007.
- [DDHR06] M. De Wulf, L. Doyen, T. Henzinger, J.-F. Raskin, *Antichains: A new algorithm for checking universality of finite automata*, Proc. CAV., LNCS **4144**, Springer, 2006, 17–30.
- [FJJ⁺12] M. Fredrikson, R. Joiner, S. Jha, T. Reps, P. Porras, H. Saïdi, and V. Yegneswaran, *Efficient runtime policy enforcement using counterexample-guided abstraction refinement*, Proc. CAV, LNCS **7358**, Springer, 2012, 548–563.
- [FKL13] O. Friedmann, F. Klaedtke, and M. Lange, *Ramsey goes visibly pushdown*, Proc. ICALP, LNCS **7966**, Springer, 2013, 224–237.
- [GN11] O. Gauwin and J. Niehren, *Streamable fragments of forward XPath*, Proc. CIAA, LNCS **6807**, Springer, 2011, 3–15.
- [GNR08] O. Gauwin, J. Niehren, and Y. Roos, *Streaming tree automata*, Information Processing Letters **109** (2008), 13–17.
- [GNT09] O. Gauwin, J. Niehren, and S. Tison, *Earliest query answering for deterministic nested word automata*, Proc. FCT, LNCS **5699**, Springer, 2009 121–132.
- [Glucose] Glucose, www.lri.fr/~simon/?page=glucose.
- [RCB08] G. Roşu, F. Chen, and T. Ball, *Runtime verification*, Proc. RV, LNCS **5289**, Springer, 2008, 51–68.
- [TO12] N.V. Tang and H. Ohsaki, *On model checking for visibly pushdown automata*, Proc. LATA, LNCS **7183**, Springer, 2012, 408–419.

A Complements to Section 3

A.1 Proof of Theorem 1

(\Rightarrow) Assume that \mathcal{A} is u -universal. Consider the set \mathcal{C} of configurations (q, σ) of \mathcal{A} such that there exists $v \in (\Sigma \cup \overline{\Sigma})^*$, $q_i \in Q_i$ and $q_f \in Q_f$ verifying $uv \in [T_\Sigma]$ and $(q_i, \epsilon) \xrightarrow{u} (q, \sigma) \xrightarrow{v} (q_f, \epsilon)$. We have $\mathcal{C} \subseteq \text{Reach}(u)$.

Let v be such that $uv \in [T_\Sigma]$. As \mathcal{A} is right-universal w.r.t. u , there exists a configuration $(q, \sigma) \in \mathcal{C}$ such that $(q, \sigma) \xrightarrow{v} (q_f, \epsilon)$ with $q_f \in Q_f$. Hence $\mathcal{C} \in \text{Safe}(u)$. As $\text{Safe}(u)$ is \subseteq -upward closed, we get $\text{Reach}(u) \in \text{Safe}(u)$.

(\Leftarrow) Assume now that $\text{Reach}(u) \in \text{Safe}(u)$, and let v be such that $uv \in [T_\Sigma]$. As $\text{Reach}(u) \in \text{Safe}(u)$, there exists $(q, \sigma) \in \text{Reach}(u)$ and $p \in Q_f$ such that $(q, \sigma) \xrightarrow{v} (p, \epsilon)$. Thus, $uv \in L(\mathcal{A})$, and \mathcal{A} is right-universal w.r.t. u .

A.2 Proof of Lemma 1

(\supseteq) Assume $\mathcal{C} \in \text{Safe}(u)$, and let v be such that $u[h]v \in [T_\Sigma]$. As h is a hedge, we have $uv \in [T_\Sigma]$. As $\mathcal{C} \in \text{Safe}(u)$, there exists $(q, \sigma) \in \mathcal{C}$ such that $(q, \sigma) \xrightarrow{v} (p, \epsilon)$ with $p \in Q_f$. So $\mathcal{C} \in \text{Safe}(u[h])$.

(\subseteq) Conversely, assume $\mathcal{C} \in \text{Safe}(u[h])$. Let v be such that $uv \in [T_\Sigma]$. We also have $u[h]v \in [T_\Sigma]$, so there exists $(q, \sigma) \in \mathcal{C}$ such that $(q, \sigma) \xrightarrow{v} (p, \epsilon)$ with $p \in Q_f$. Thus $\mathcal{C} \in \text{Safe}(u)$.

The proof is the same for $L\text{Safe}(u[h]) = L\text{Safe}(u)$, except that we only consider v of the form $\overline{a}v'$.

A.3 Proof of Proposition 1

Before giving the proof of Proposition 1, we introduce the *Post* operator and Proposition 2. From $u \in (\Sigma \cup \overline{\Sigma})^*$ and a set of configurations $\mathcal{C} \subseteq Q \times \Gamma^*$, we define $\text{Post}_u(\mathcal{C})$ as the set

$$\text{Post}_u(\mathcal{C}) = \{(q', \sigma') \in Q \times \Gamma^* \mid \exists (q, \sigma) \in \mathcal{C}, (q, \sigma) \xrightarrow{u} (q', \sigma')\}.$$

Proposition 2. *Let $ua \in \text{Pref}(T_\Sigma)$ with $a \in \Sigma$.*

$$\mathcal{C} \in L\text{Safe}(ua) \iff \text{Post}_{\overline{a}}(\mathcal{C}) \in \text{Safe}(u) \quad (15)$$

$$\mathcal{C} \in \text{Safe}(ua) \iff \forall h \in H_\Sigma, \text{Post}_{[h]}(\mathcal{C}) \in L\text{Safe}(ua) \quad (16)$$

Proof. (15, \Rightarrow) Let $\mathcal{C} \in L\text{Safe}(ua)$ and $\mathcal{C}' = \text{Post}_{\overline{a}}(\mathcal{C})$. Let us show that $\mathcal{C}' \in \text{Safe}(u)$. By Lemma 1, it is sufficient to prove that $\mathcal{C}' \in \text{Safe}(ua\overline{a})$. Let v such that $ua\overline{a}v \in [T_\Sigma]$. As $\mathcal{C} \in L\text{Safe}(ua)$ and $\overline{a}v$ starts with $\overline{a} \in \overline{\Sigma}$, there exists $(q, \sigma) \in \mathcal{C}$ and (q', σ') such that $(q, \sigma) \xrightarrow{\overline{a}} (q', \sigma') \xrightarrow{v} (p, \epsilon)$ for some $p \in Q_f$. By definition of $\text{Post}_{\overline{a}}(\mathcal{C})$ we have $(q', \sigma') \in \mathcal{C}'$ and thus $\mathcal{C}' \in \text{Safe}(ua\overline{a})$.

(15, \Leftarrow) For the converse, let $\mathcal{C}' = \text{Post}_{\overline{a}}(\mathcal{C}) \in \text{Safe}(u) = \text{Safe}(ua\overline{a})$. Let us show that $\mathcal{C} \in L\text{Safe}(ua)$. Let v be such that $uav \in [T_\Sigma]$ and $v = \overline{b}v'$. We necessarily have $\overline{a} = \overline{b}$. As $\mathcal{C}' \in \text{Safe}(ua\overline{a})$, there exists $(q', \sigma') \in \mathcal{C}'$ such

that $(q', \sigma') \xrightarrow{v'} (p, \epsilon)$ with $p \in Q_f$. By definition of $Post_{\bar{a}}(\mathcal{C})$, there also exists $(q, \sigma) \in \mathcal{C}$ such that $(q, \sigma) \xrightarrow{\bar{a}} (q', \sigma')$ and thus $(q, \sigma) \xrightarrow{v=\bar{a}v'} (p, \epsilon)$ with $p \in Q_f$.

(16, \Rightarrow) Let $\mathcal{C} \in Safe(ua)$ and $h \in H_\Sigma$. Let us show that $\mathcal{C}' = Post_{[h]}(\mathcal{C})$ is in $LSafe(ua)$. Let v such that $uav \in [T_\Sigma]$ and $v = \bar{b}v'$. We must have $\bar{a} = \bar{b}$. We also have $ua[h]\bar{a}v' \in [T_\Sigma]$. As $\mathcal{C} \in Safe(ua)$, there exists $(q, \sigma) \in \mathcal{C}$ such that $(q, \sigma) \xrightarrow{[h]} (q', \sigma') \xrightarrow{v=\bar{a}v'} (p, \epsilon)$ with $p \in Q_f$. By definition of $Post_{[h]}(\mathcal{C})$, $(q', \sigma') \in \mathcal{C}'$.

(16, \Leftarrow) Let us assume that for every hedge $h \in H_\Sigma$, $Post_{[h]}(\mathcal{C}) \in LSafe(ua)$. Let us show that $\mathcal{C} \in Safe(ua)$. Let v be such that $uav \in [T_\Sigma]$. Then we have $uav = ua[h]\bar{a}v'$ for some $h \in H_\Sigma$. As $\mathcal{C}' = Post_{[h]}(\mathcal{C}) \in LSafe(ua)$, $\bar{a}v'$ starts with $\bar{a} \in \bar{\Sigma}$ and $ua\bar{a}v' \in [T_\Sigma]$, there exists $(q', \sigma') \in \mathcal{C}'$ such that $(q', \sigma') \xrightarrow{\bar{a}v'} (p, \epsilon)$ for some $p \in Q_f$. Hence, by definition of $Post_{[h]}(\mathcal{C})$, there also exists $(q, \sigma) \in \mathcal{C}$ such that $(q, \sigma) \xrightarrow{[h]} (q', \sigma') \xrightarrow{\bar{a}v'} (p, \epsilon)$ with $p \in Q_f$. \square

Predecessors closely relate to the *Post* operator, in the following sense.

Lemma 4. *\mathcal{C} is an \bar{a} -predecessor of $Post_{\bar{a}}(\mathcal{C})$. If $\mathcal{C} \in Pred_{\bar{a}}(\mathcal{C}')$, then $\mathcal{C}' \subseteq Post_{\bar{a}}(\mathcal{C})$. Both properties also hold for $Post_{[h]}(\mathcal{C})$.*

The proof of Proposition 1 is then the following one.

Proof (of Proposition 1). (3, \Rightarrow) Let $\mathcal{C} \in LSafe(ua)$. Then by Proposition 2, $Post_{\bar{a}}(\mathcal{C}) \in Safe(u)$. Moreover, \mathcal{C} is an \bar{a} -predecessor of $Post_{\bar{a}}(\mathcal{C})$ by Lemma 4.

(3, \Leftarrow) Let \mathcal{C} be an \bar{a} -predecessor of \mathcal{C}' , with $\mathcal{C}' \in Safe(u)$. By Lemma 4, $\mathcal{C}' \subseteq Post_{\bar{a}}(\mathcal{C})$. As $Safe(u)$ is \subseteq -upward closed, we also have $Post_{\bar{a}}(\mathcal{C}) \in Safe(u)$, so $\mathcal{C} \in LSafe(ua)$ by Proposition 2.

(4) Same proofs, except that \bar{a} has to be replaced by h , for all $h \in H_\Sigma$. \square

B Complements to Section 4

B.1 Proof of Lemma 2

Let us recall the definition of h -predecessor: \mathcal{C} is a h -predecessor of \mathcal{C}' if $\forall (q', \sigma) \in \mathcal{C}'$, $\exists (q, \sigma) \in \mathcal{C}$, $(q, \sigma) \xrightarrow{h} (q', \sigma)$. Hence if $h \sim h'$, then \mathcal{C} is a h -predecessor of \mathcal{C}' iff \mathcal{C} is a h' -predecessor of \mathcal{C}' .

B.2 Algorithm computing H

We propose an algorithm for computing such a set H from a VPA \mathcal{A} . Algorithm 3 is based on the definition of hedges, adapted to their associated functions:

- ϵ is the empty hedge, and $rel_\epsilon(q) = \{q\}$ for every $q \in Q$. We write this function id_Q .
- if h_1, h_2 are two hedges, then h_1h_2 is a hedge, and $rel_{h_1h_2} = rel_{h_2} \circ rel_{h_1}$.

- if h is a hedge and $a \in \Sigma$, then $ah\bar{a}$ is a hedge, and $\text{rel}_{ah\bar{a}}(q)$ is the set of states q' such that there exists $\gamma \in \Gamma$ verifying:

$$(q, \epsilon) \xrightarrow{a} (p, \gamma) \quad \text{and} \quad (p', \gamma) \xrightarrow{\bar{a}} (q', \epsilon) \quad \text{with } p' \in \text{rel}_h(p).$$

Algorithm 3 uses the variables *ToProcess* and *Functions* with the following meaning. *Functions* contains initially the identity relation id_Q ; at the end of the computation, it contains all functions rel_h , for $h \in H_\Sigma$. *ToProcess* contains all the newly constructed functions, and these functions are used to create other new functions as described in the previous definition by induction.

Algorithm 3 Computing all functions rel_h , for $h \in H_\Sigma$.

```

function HEDGEFUNCTIONS( $\mathcal{A}$ )
  Functions  $\leftarrow \{\text{id}_Q\}$ 
  ToProcess  $\leftarrow \{\text{id}_Q\}$ 
  while ToProcess  $\neq \emptyset$  do
    fct  $\leftarrow \text{POP}(\textit{ToProcess})$ 
    NewFunctions  $\leftarrow \emptyset$ 
    for  $f \in \textit{Functions}$  do
      NewFunctions  $\leftarrow \textit{NewFunctions} \cup \{f \circ \textit{fct}, \textit{fct} \circ f\}$ 
    end for
    for  $a \in \Sigma$  do
       $f \leftarrow f_\emptyset$  //  $f_\emptyset$  maps every  $q \in Q$  to  $\emptyset$ 
      for  $q \xrightarrow{a:\gamma} p \in \Delta$  and  $p' \xrightarrow{\bar{a}:\gamma} q' \in \Delta$  with  $p' \in \textit{fct}(p)$  do
         $f(q) \leftarrow f(q) \cup \{q'\}$ 
      end for
      NewFunctions  $\leftarrow \textit{NewFunctions} \cup \{f\}$ 
    end for
    ToProcess  $\leftarrow \textit{ToProcess} \cup (\textit{NewFunctions} \setminus \textit{Functions})$ 
    Functions  $\leftarrow \textit{Functions} \cup \textit{NewFunctions}$ 
  end while
  return Functions
end function

```

Proposition 3. *Algorithm 3 computes the set $\{\text{rel}_h \mid h \in H_\Sigma\}$.*

Proof. Let *Functions* be the set computed by Algorithm 3. Clearly, $\text{Functions} \subseteq \{\text{rel}_h \mid h \in H_\Sigma\}$. Assume for contradiction that there exists $r = \text{rel}_h$ with $h \in H_\Sigma$ such that $r \notin \text{Functions}$. Clearly, $r \neq \text{id}_Q$, and we can suppose wlog that either $r = r'_2 \circ r'_1$ with $r'_1, r'_2 \in \text{Functions} \setminus \{\text{id}_Q\}$, or there exists $r' \in \text{Functions}$ such that for all q , $r(q)$ is the set of q' with $q \xrightarrow{a:\gamma} p \in \Delta$, $p' \xrightarrow{\bar{a}:\gamma} q' \in \Delta$ and $p' \in r'(p)$. Consider the first case. When they have been constructed by Algorithm 3, both r'_1 and r'_2 have been added to *ToProcess* and to *Functions*. After the last element (among r'_1 and r'_2) is popped from *ToProcess*, then $r = r'_2 \circ r'_1$ is built during the loop on $f \in \text{Functions}$, which leads to a contradiction. We also have a contradiction in the second case by considering the loop on $a \in \Sigma$. \square

B.3 Antichains $[Safe(u)]$ and $[LSafe(u)]$ are Finite

Proposition 4. $[Safe(u)]$ and $[LSafe(u)]$ are finite and only contain finite sets of configurations.

Proof. We begin with the following observation. Let v be such that $[uv] \in T_\Sigma$ and $(q, \sigma) \xrightarrow{v} (p, \epsilon)$ with $p \in Q_f$. Let u' (resp. v') the word obtained from u (resp. v) by removing all factors that are linearizations of hedges. Then $|u'| = |v'|$ and $|u'| = |\sigma|$.

Let $\mathcal{C} \in Safe(u)$. Then by definition

$$\forall v, uv \in [T_\Sigma] \implies \exists (q, \sigma) \in \mathcal{C}, (q, \sigma) \xrightarrow{v} (p, \epsilon) \text{ with } p \in Q_f.$$

If \mathcal{C} is \subseteq -minimal, then every $(q, \sigma) \in \mathcal{C}$ is used for at least one v in the previous definition. Now by the previous observation, each such (q, σ) belongs to $Q \times \Gamma^{|u'|}$. Hence $\mathcal{C} \subseteq Q \times \Gamma^{|u'|}$, and thus both \mathcal{C} and $[Safe(u)]$ are finite.

The same arguments hold for proving that $[LSafe(u)]$ is finite and contains only finite sets of configurations. \square

B.4 Proof of Implication (6)

Let $\mathcal{C} \in [LSafe(ua)]$ and let $\mathcal{C}' = Post_{\bar{a}}(\mathcal{C})$. We know from Proposition 2 that $\mathcal{C}' \in Safe(u)$. Let $\mathcal{C}'_0 \subseteq \mathcal{C}'$ such that $\mathcal{C}'_0 \in [Safe(u)]$. From the definition of \mathcal{C}' we get:

$$\forall c' \in \mathcal{C}', \exists c \in \mathcal{C}, c \xrightarrow{\bar{a}} c'$$

We build \mathcal{C}_0 from these $c \in \mathcal{C}$ but for $c' \in \mathcal{C}'_0$:

$$\mathcal{C}_0 = \{c \in \mathcal{C} \mid \exists c' \in \mathcal{C}'_0, c \xrightarrow{\bar{a}} c'\}$$

Figure 5 illustrates the construction. The set \mathcal{C}_0 is an \bar{a} -predecessor of \mathcal{C}'_0 , so us-

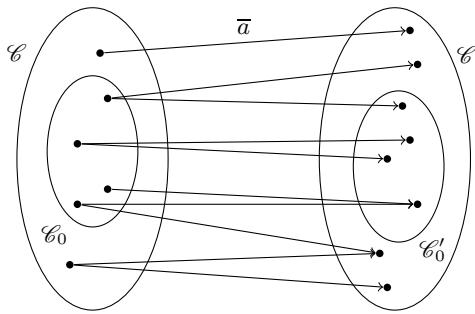


Fig. 5: Construction of \mathcal{C}_0

ing Proposition 1, we get $\mathcal{C}_0 \in LSafe(ua)$. Furthermore, $\mathcal{C}_0 \subseteq \mathcal{C} \in [LSafe(ua)]$, so $\mathcal{C}_0 = \mathcal{C}$, and \mathcal{C} is obtained as an \bar{a} -predecessor of $\mathcal{C}'_0 \in [Safe(u)]$.

B.5 Proof of Implication (8)

The proof is in the same vein as for Implication (6). Let $\mathcal{C} \in \llbracket \text{Safe}(ua) \rrbracket$, and $h \in H_\Sigma$. Let $\mathcal{C}'_h = \text{Post}_{[h]}(\mathcal{C})$. By Proposition 2, $\mathcal{C}'_h \in \text{LSafe}(ua)$. Let $\mathcal{C}''_h \subseteq \mathcal{C}'_h$ such that $\mathcal{C}''_h \in \llbracket \text{LSafe}(ua) \rrbracket$. We know that $\forall c' \in \mathcal{C}'_h, \exists c \in \mathcal{C}$ such that $c \xrightarrow{[h]} c'$. We define $\mathcal{C}_h = \{c \in \mathcal{C} \mid \exists c' \in \mathcal{C}''_h, c \xrightarrow{[h]} c'\}$. For every $h \in H_\Sigma$, \mathcal{C}_h is a h -predecessor of $\mathcal{C}''_h \in \text{LSafe}(ua)$. Consider $\mathcal{C}_\cup = \bigcup_{h \in H_\Sigma} \mathcal{C}_h$, then \mathcal{C}_\cup is also a h -predecessor of \mathcal{C}''_h . Using Proposition 1, we have $\mathcal{C}_\cup \in \text{Safe}(ua)$. As $\mathcal{C}_\cup \subseteq \mathcal{C}$ and $\mathcal{C} \in \llbracket \text{Safe}(ua) \rrbracket$, we also have that $\mathcal{C}_\cup = \mathcal{C}$. Hence \mathcal{C} verifies that $\forall h \in H_\Sigma, \exists \mathcal{C}'' \in \llbracket \text{LSafe}(ua) \rrbracket$ such that \mathcal{C} is a h -predecessor of \mathcal{C}'' .

C Complements to Section 5

C.1 Proof of Equation (11)

Let S denote the set

$$\left\{ \mathcal{C} \mid \mathcal{C} = \bigcup_{h \in H} \mathcal{C}_h \text{ with } \mathcal{C}_h \text{ a minimal } h\text{-predecessor of } \mathcal{C}' \in \llbracket \text{LSafe}(ua) \rrbracket \right\}$$

Let $\mathcal{C} \in S$. We have: $\mathcal{C} = \underbrace{\mathcal{C}_{h_1} \cup \dots \cup \mathcal{C}_{h_k}}_{h_i \in [H]} \cup \underbrace{\mathcal{C}_{h'_1} \cup \dots \cup \mathcal{C}_{h'_n}}_{h'_i \in H \setminus [H]}$. Let us show that

$\mathcal{C}_{h_1} \cup \dots \cup \mathcal{C}_{h_k} \cup \mathcal{C}_{h'_1} \cup \dots \cup \mathcal{C}_{h'_{n-1}} \in S$ ($\mathcal{C}_{h'_n}$ has been erased). By induction, this will prove that $\mathcal{C}_{h_1} \cup \dots \cup \mathcal{C}_{h_k} \in S$. We have $h'_n \in H \setminus [H]$, so there exists $h_i \in [H]$ such that $h_i \leq h'_n$. As \mathcal{C}_{h_i} is a minimal h_i -predecessor of an element \mathcal{C}' in $\llbracket \text{LSafe}(ua) \rrbracket$, it follows from (10) that \mathcal{C}_{h_i} is also a minimal h'_n -predecessor of \mathcal{C}' . So $\mathcal{C}_{h_1} \cup \dots \cup \mathcal{C}_{h_k} \cup \mathcal{C}_{h'_1} \cup \dots \cup \mathcal{C}_{h'_{n-1}} \cup \mathcal{C}_{h_i} \in S$.

C.2 Proof of Equation (13)

(\supseteq) Let $\mathcal{C} \in \llbracket (A \cap B) \cup (A \setminus B \sqcup B \setminus A) \rrbracket$. Then $\mathcal{C} \in A \sqcup B$. For contradiction, let us assume that there exists $\mathcal{C}' \subsetneq \mathcal{C}$ such that $\mathcal{C}' \in \llbracket A \sqcup B \rrbracket$. If $\mathcal{C}' \in A \cap B$, then $\mathcal{C}' \in (A \cap B) \cup (A \setminus B \sqcup B \setminus A)$, which contradicts \mathcal{C} . So $\mathcal{C}' \notin A \cap B$, and assume wlog that $\mathcal{C}' = a \cup b$ with $a \in A \setminus B$ and $b \in B$ (recall that a, b are non empty sets by definition of \sqcup operator). If $b \in A$, then $b \in A \cap B \subseteq A \sqcup B$ and $b \subsetneq \mathcal{C}'$, but this contradicts \mathcal{C}' . If $b \notin A$, then $\mathcal{C}' \in A \setminus B \sqcup B \setminus A$, so $\mathcal{C}' \in A \cap B \cup (A \setminus B \sqcup B \setminus A)$, and $\mathcal{C}' \subsetneq \mathcal{C}$, which contradicts \mathcal{C} .

(\subseteq) Let $\mathcal{C} \in \llbracket A \sqcup B \rrbracket$. Let us first show that $\mathcal{C} \in (A \cap B) \cup (A \setminus B \sqcup B \setminus A)$. If $\mathcal{C} \in A \cap B$ this is direct. Otherwise $\mathcal{C} = a \cup b$ with $a \in A \setminus B$ and $b \in B$ (the other case is symmetric). If $b \in A$ then $b \in A \cap B \subseteq A \sqcup B$ and $b \subsetneq \mathcal{C}$, which contradicts the definition of \mathcal{C} . So $b \in B \setminus A$, and $\mathcal{C} \in A \setminus B \sqcup B \setminus A$. Secondly, assume for contradiction that there exists $\mathcal{C}' \subsetneq \mathcal{C}$ such that $\mathcal{C}' \in \llbracket (A \cap B) \cup (A \setminus B \sqcup B \setminus A) \rrbracket$. Then, according to (\supseteq), $\mathcal{C}' \in \llbracket A \sqcup B \rrbracket$, which contradicts the definition of \mathcal{C} .

We get the next corollary of Equation (13):

Corollary 2. *If $A \subseteq B$, then $\llbracket A \sqcup B \rrbracket = \llbracket A \rrbracket$.*

C.3 *Safe(a)* and Operator \sqcup

The \sqcup operator also simplifies the definition of $\lfloor \text{Safe}(a) \rfloor$.

Proposition 5. $\lfloor \text{Safe}(a) \rfloor = \left\lfloor \bigsqcup_{h \in [H]} A_h \right\rfloor$ with

$$A_h = \left\{ \{(q, \sigma)\} \mid \exists q_f \in Q_f : (q, \sigma) \xrightarrow{h\bar{a}} (q_f, \epsilon) \right\}.$$

Proof. We consider $\lfloor \text{Safe}(a) \rfloor$ as defined in Equation (12). Let S denote the set $\left\{ \mathcal{C} \mid \forall h \in [H], \exists q_f \in Q_f, \exists (q, \sigma) \in \mathcal{C} : (q, \sigma) \xrightarrow{h\bar{a}} (q_f, \epsilon) \right\}$.

1. Every element of $\bigsqcup_{h \in [H]} A_h$ belongs to S . Thus $\left\lfloor \bigsqcup_{h \in [H]} A_h \right\rfloor \subseteq S$.
2. Let us show that for each \mathcal{C} in S , there exists $\mathcal{C}' \in \bigsqcup_{h \in [H]} A_h$ such that $\mathcal{C}' \subseteq \mathcal{C}$. Let $\mathcal{C} \in S$. By definition of S , for all $h \in [H]$ there exists $(q_h, \sigma_h) \in \mathcal{C}$ and $q_f \in Q_f$ such that $(q_h, \sigma_h) \xrightarrow{h\bar{a}} (q_f, \epsilon)$. Let $\mathcal{C}' = \{(q_h, \sigma_h) \mid h \in [H]\}$. Then $\mathcal{C}' \subseteq \mathcal{C}$ and $\mathcal{C}' \in \bigsqcup_{h \in [H]} A_h$ because $\{(q_h, \sigma_h)\} \in A_h, \forall h$.
3. Assume that there exists $\mathcal{C}_* \in \left\lfloor \bigsqcup_{h \in [H]} A_h \right\rfloor \setminus [S]$. By 1., there exists \mathcal{C} in $[S]$ such that $\mathcal{C} \subsetneq \mathcal{C}_*$; and by 2., there exists $\mathcal{C}' \in \bigsqcup_{h \in [H]} A_h$ such that $\mathcal{C}' \subseteq \mathcal{C} \subsetneq \mathcal{C}_*$ in contradiction with the definition of \mathcal{C}_* . Therefore $\left\lfloor \bigsqcup_{h \in [H]} A_h \right\rfloor \subseteq [S]$.
4. Let $\mathcal{C} \in [S]$. By 2., there exists $\mathcal{C}' \in \left\lfloor \bigsqcup_{h \in [H]} A_h \right\rfloor$ such that $\mathcal{C}' \subseteq \mathcal{C}$. By 3., it follows that $\mathcal{C} = \mathcal{C}'$ and thus $[S] \subseteq \left\lfloor \bigsqcup_{h \in [H]} A_h \right\rfloor$.

In conclusion, $\lfloor \text{Safe}(a) \rfloor = [S] = \left\lfloor \bigsqcup_{h \in [H]} A_h \right\rfloor$ □

C.4 Step 2 with SAT Solvers

The computation of the set $\lfloor \text{Safe}(ua) \rfloor$ from $\lfloor \text{LSafe}(ua) \rfloor$ can be done with a SAT solver as described in Equation (11) :

$$\lfloor \text{Safe}(ua) \rfloor = \left\lfloor \left\{ \mathcal{C} \mid \mathcal{C} = \bigcup_{h \in [H]} \mathcal{C}_h \text{ with } \mathcal{C}_h \text{ a } h\text{-predecessor of } \mathcal{C}' \in \lfloor \text{LSafe}(ua) \rfloor \right\} \right\rfloor.$$

The method is the same as for Step 1, with the only difference that we use the next formula instead of formula $\varphi_{\bar{a}}$:

$$\bigwedge_{h \in [H]} \bigvee_{\mathcal{C}' \in \lfloor \text{LSafe}(ua) \rfloor} \bigwedge_{c' \in \mathcal{C}'} \bigvee_{c \xrightarrow{h} c'} x_c,$$

D Complements to Section 6

We provide in this section additional experimental results.

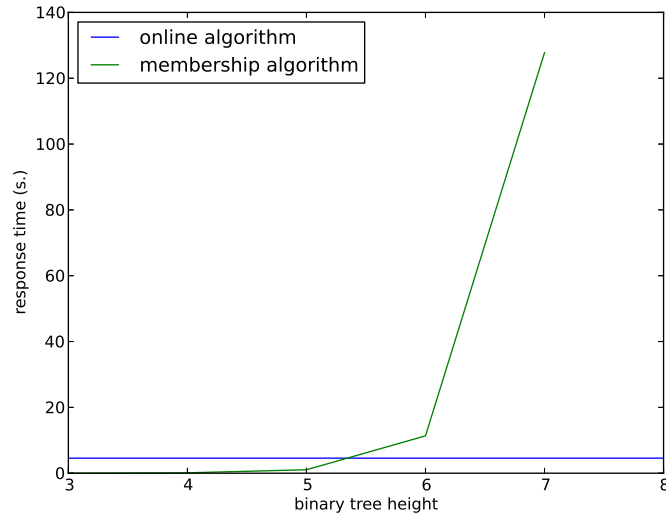


Fig. 6: Execution time of the membership algorithm for a binary tree t_0 of increasing height.

D.1 Randomly Generated VPAs

Figure 6 indicates the execution time of the membership algorithm for a random complete binary tree t_0 of increasing height.

D.2 VPAs resulting from XPath translation

Figure 7 indicates the computation time of our online algorithm for 15 random trees, with a height bounded by 10, and a number of children by node bounded by 15.

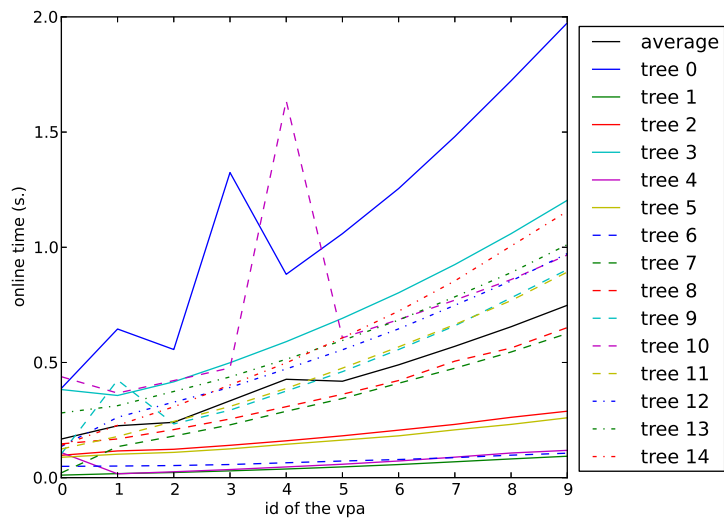


Fig. 7: Computation time of the online algorithm for each random tree.