

Regular Paper

An On-The-Fly Algorithm for Conditional Weighted Pushdown Systems

HUA VY LE THANH^{1,a)} XIN LI^{2,b)}

Received: January 14, 2014, Accepted: April 18, 2014

Abstract: Pushdown systems (PDSs) are well-understood as abstract models of recursive sequential programs, and weighted pushdown systems (WPDSs) are a general framework for solving certain meet-over-all-path problems in program analysis. Conditional WPDSs (CWPDSs) further extend WPDSs to enhance the expressiveness of WPDSs, in which each transition is guarded by a regular language over the stack that specifies conditions under which a transition rule can be applied. CWPDSs or its instance are shown to have wide applications in analysis of objected-oriented programs, access rights analysis, etc. Model checking CWPDSs was shown to be reduced to model checking WPDSs, and an offline algorithm was given that translates CWPDSs to WPDSs by synchronizing the underlying PDS and finite state automata accepting regular conditions. The translation, however, can cause an exponential blow-up of the system. This paper presents an on-the-fly model checking algorithm for CWPDSs that synchronizes the computing machineries on-demand while computing post-images of regular configurations. We developed an on-the-fly model checker for CWPDSs and apply it to models generated from the reachability analysis of the HTML5 parser specification. Our preliminary experiments show that, the on-the-fly algorithm drastically outperforms the offline algorithm regarding both practical space and time efficiency.

Keywords: Conditional Weighted Pushdown System, model checking

1. Introduction

Pushdown systems (PDSs) are well understood as abstract models of sequential programs with recursive procedures. By encoding programs as pushdown systems, procedure calls and returns are guaranteed to be correctly paired. Weighted pushdown systems (WPDSs) [4] extend PDSs, by associating a value with each transition, and are a general framework for solving certain meet-over-all-path problems in program analysis. Program analysis yielded by WPDSs are context-sensitive in terms of valid paths. Li et al. [1] further extended WPDSs to Conditional Weighted Pushdown Systems (CWPDSs), in which each transition is guarded by a regular language that specifies conditions under which the transition can be applied. The extension is motivated by the observation that, WPDSs are not precise enough to model objected-oriented programs like Java in program analysis.

CWPDSs are more expressive than WPDSs, and have wider applications. A points-to analysis algorithm was designed in the framework of CWPDSs, which enjoys context-sensitivity beyond valid paths regarding heap abstraction, call graph construction, and heap access [1]. Besides modelling objected-oriented program features, it can yield algorithms for generating stack-based access control policies for programs that employ stack-based ac-

cess control mechanism, such as Java or and Microsoft CLR (common language runtime) [7]. Minamide et al. [2] adopted conditional pushdown systems (CPDSs), an instance of CWPDSs, to model the HTML5 parser specification, and performed reachability analysis over it. They proposed a new algorithm for reachability analysis of CPDSs and successfully found several compatibility problems of web browsers with HTML5 specifications.

It was shown that the model checking problem on CWPDSs can be reduced to that on WPDSs, and an offline algorithm was given, which translates CWPDSs to WPDSs by synchronizing the underlying pushdown system and finite state automata built for regular conditions [1]. The translation, however, can cause an exponential blow-up of the system, and prohibits the offline algorithm from practical applications. The original offline algorithm for CWPDSs generated the entire state space, whereas only part of the space needs exploring to answer a model checking enquiry.

This paper makes the following contributions:

- We present an on-the-fly model checking algorithm for CWPDSs, by interleaving saturation procedures for computing post-images with regular condition checking, and synchronizing the underlying pushdown systems and finite state automata accepting regular conditions on-demand. We also provide concrete algorithms for those steps of model checking CWPDSs that are only defined in the original paper.
- We implemented the on-the-fly model checking algorithm in terms of computing post-images, and evaluated the algorithm with one of the large examples analysed in reachability analysis of HTML5 parser specification. Our preliminary

¹ University of Science - Ho Chi Minh City, Vietnam.

² Corresponding author. The University of Tokyo, Bunkyo, Tokyo 113-8654, Japan.

a) thanhvy@gmail.com

b) li-xin@kb.is.s.u-tokyo.ac.jp

experiments showed that the on-the-fly algorithm drastically outperforms the offline algorithm regarding both space and time efficiency in practice.

Related Work

Esparza et al introduced solutions for model checking LTL on pushdown systems with regular valuation. As an application, they introduced a formal model called pushdown systems with checkpoints, and showed that the reachability problem of the model is EXPTIME-complete. The underlying CPDSs of CWPDSs can be regarded as an instance of the model, and also shares the EXPTIME-hardness for reachability checking. So our on-the-fly algorithm for CWPDSs is designed to hopefully be practical by exploring the state space relevant to model checking enquiries.

Minamide et al. adopted CPDSs as a formal model for reachability analysis of HTML5 parser specification and then checked compatibility of web browsers with specifications. The authors proposed a new algorithm for computing regular pre-images of CPDSs, by extending \mathcal{P} -automata with regular lookahead. Though the algorithm also has exponential complexity, it aimed to avoid the exponential blow-up of the system state space due to the translation. Preliminary experimental results in Section 4 indicate that our on-the-fly algorithm is also competitive in avoiding the exponential blow-up before reachability analysis.

Organization The remainder of the paper is organized as follows. Section 2 recalls CWPDSs and the off-line model checking algorithm. Section 3 presents an on-the-fly algorithm for CWPDSs, and an example that illustrates the procedure. Our empirical study is given in Section 4, which compares performance of the on-the-fly algorithm and the off-line algorithm. Section 5 concludes the paper.

2. Preliminary

2.1 Weighted Pushdown Systems

Definition 2.1 A **pushdown system** \mathcal{P} is $(Q, \Gamma, \Delta, q_0, \omega_0)$ where Q is a finite set of control locations, Γ is a finite stack alphabet, $\Delta \subseteq Q \times \Gamma \times Q \times \Gamma^*$ is a finite set of transitions, $q_0 \in Q$ and $\omega_0 \in \Gamma^*$ is the initial control location and stack contents, respectively. A transition $(p, \gamma, q, \omega) \in \Delta$ is written as $\langle p, \gamma \rangle \hookrightarrow \langle q, \omega \rangle$. A **configuration** is a pair $\langle q, \omega \rangle$ with $q \in Q$ and $\omega \in \Gamma^*$. A computation relation \Rightarrow is defined on configurations, such that $\langle p, \gamma \omega' \rangle \Rightarrow \langle q, \omega \omega' \rangle$ for all $\omega' \in \Gamma^*$ if $\langle p, \gamma \rangle \hookrightarrow \langle q, \omega \rangle$, and the reflexive and transitive closure of \Rightarrow is denoted by \Rightarrow^* . Given a set of configurations C , we define $pre^*(C) = \{c' \mid \exists c \in C : c' \Rightarrow^* c\}$ and $post^*(C) = \{c' \mid \exists c \in C : c \Rightarrow^* c'\}$ to be the (possibly infinite) set of pre- and post- images of C , respectively.

A pushdown system is a variant of pushdown automata without input alphabet. It is not used as language acceptors but a model of system behaviors.

Definition 2.2 A **bounded idempotent semiring** \mathcal{S} is $(D, \oplus, \otimes, \bar{0}, \bar{1})$, where $\bar{0}, \bar{1} \in D$, and

- (D, \oplus) is a commutative monoid with $\bar{0}$ as its unit element, and \oplus is idempotent, i.e., $a \oplus a = a$ for all $a \in D$;
- (D, \otimes) is a monoid with $\bar{1}$ as the unit element;
- \otimes distributes over \oplus , i.e., for all $a, b, c \in D$, $a \otimes (b \oplus c) = (a \otimes b) \oplus (a \otimes c)$ and $(b \oplus c) \otimes a = (b \otimes a) \oplus (c \otimes a)$;
- for all $a \in D$, $a \otimes \bar{0} = \bar{0} \otimes a = \bar{0}$;

- A partial ordering \sqsubseteq is defined on D such that $a \sqsubseteq b$ iff $a \oplus b = a$ for all $a, b \in D$, and there are no infinite descending chains in D .

Definition 2.3 A **weighted pushdown system** \mathcal{W} is $(\mathcal{P}, \mathcal{S}, f)$, where $\mathcal{P} = (Q, \Gamma, \Delta, p_0, \omega_0)$ is a pushdown system, $\mathcal{S} = (D, \oplus, \otimes, \bar{0}, \bar{1})$ is a bounded idempotent semiring, and $f : \Delta \rightarrow D$ is a weight assignment function.

By Definition 2.2, we have that $\bar{0}$ is the greatest element.

Definition 2.4 Given a weighted pushdown system $\mathcal{W} = (\mathcal{P}, \mathcal{S}, f)$ where $\mathcal{P} = (Q, \Gamma, \Delta, p_0, \omega_0)$. Let $\sigma = [r_1, \dots, r_k] \in \Delta^*$ denote a sequence of pushdown transition rules, and define $val(\sigma) = f(r_1) \otimes \dots \otimes f(r_k)$. For any $c, c' \in Q \times \Gamma^*$, we denote by $path(c, c')$ the set of sequences of transition rules that transform configurations from c into c' . Given regular sets of configurations $S, T \subseteq Q \times \Gamma^*$, the **meet-over-all-valid-path (MOVP)** problem computes

$$MOVP(S, T) = \oplus \{val(\sigma) \mid \sigma \in path(s, t), s \in S, t \in T\}$$

WPDSs [4] solve the MOVP problem in program analysis, of which data domains comply with the bounded idempotent semiring. When applying WPDSs to program analysis, PDSs typically model (recursive) control flows of the program, weights encode transfer functions, \otimes corresponds to the reverse of function composition, and \oplus joins data flows.

2.2 Computing Post Images

Definition 2.5 Given a PDS $\mathcal{P} = (Q, \Gamma, \Delta, q_0, \omega_0)$. A \mathcal{P} -automaton $\mathcal{A} = (Q', \Gamma, \rightarrow, Q, F)$ is a non-deterministic finite automaton, where Q' is the finite set of states, Γ is the input alphabet, $\rightarrow \subseteq Q' \times \Gamma \times Q'$ is the set of transitions, and $Q \subseteq Q'$ and $F \subseteq Q'$ are the set of initial and final states, respectively. We define $\rightarrow^* \subseteq Q' \times \Gamma^* \times Q'$ as the smallest relation such that

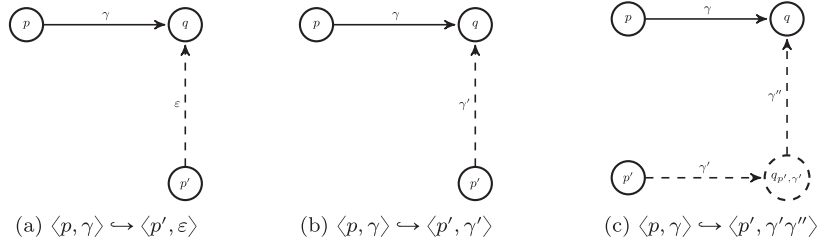
- $p \xrightarrow{\epsilon}^* p$ for any $p \in Q'$;
- $p \xrightarrow{\gamma}^* p'$ if $(p, \gamma, p') \in \rightarrow$;
- $p \xrightarrow{\omega\gamma}^* p'$ if $p \xrightarrow{\omega}^* p''$ and $p'' \xrightarrow{\gamma}^* p'$ for some $p'' \in Q'$.

A configuration $\langle p, \omega \rangle$ is accepted by \mathcal{A} if $p \xrightarrow{\omega}^* q$ for some $q \in F$. A set C of configurations is accepted by \mathcal{A} if each $c \in C$ is accepted by \mathcal{A} . A set C of configurations is regular if it is accepted by some \mathcal{P} -automaton, denoted by \mathcal{A}_C .

A key property of pushdown systems is that, a regular set of configurations is closed under forward and backward reachability. Therefore infinite configurations can be represented by finite means of \mathcal{P} -automata, and pre- and post-images can be computed by saturating the automaton.

We illustrate in **Fig.1** the saturation rules for computing $post^*(C)$ of a regular set C of configurations. The saturation procedure takes as input the \mathcal{P} -automaton \mathcal{A}_C that recognizes C , and augments \mathcal{A}_C with new edges and states by the saturation rules until convergence, and the stabilised automaton accepts $post^*(C)$. In the figure, solid edges and states reside in the current automaton, and dashed edges and states are newly added by saturation rules.

The MOVP problem tackled by WPDSs can be solved by saturating a *weighted* extension of \mathcal{P} -automata when computing pre- or post-images. Consider computing post-images as given in Algorithm 2 and ignore those underlined parts. Initially, all transi-


 Fig. 1 Saturation Rules for Computing $post^*(C)$.

tions in \mathcal{A}_C are labelled with $\bar{1}$. The saturation procedure computes or updates the value for each transition when applying the saturation rule at line 18, 21, 27 and 30, and outputs a weighted automaton $\mathcal{A}_{post^*(C)}$ of which each transition is associated with a value from D . The final result computed by the MOVP problem can be read-off from $\mathcal{A}_{post^*(C)}$. Note that the time complexity is increased from the unweighted case by a factor no more than the maximal-length of descending chains in D [4].

2.3 Conditional Weighted Pushdown Systems

Definition 2.6 A conditional pushdown system (CPDS) \mathcal{P}_c is a 6-tuple $(Q, \Gamma, C, \Delta_c, q_0, \omega_0)$, where Q is a finite set of control locations, Γ is a finite stack alphabet, C is a finite set of regular languages over Γ , $\Delta_c \subseteq Q \times \Gamma \times C \times Q \times \Gamma^*$ is the set of transition rules, and $q_0 \in Q$ and $\omega_0 \in \Gamma^*$ are the initial control location and stack contents, respectively. We write $\langle p, \gamma \rangle \xrightarrow{L} \langle q, \omega \rangle$ if $(p, \gamma, L, q, \omega) \in \Delta_c$. A computation relation \Rightarrow_c on configurations is defined such that $\langle p, \gamma \omega' \rangle \Rightarrow_c \langle q, \omega \omega' \rangle$ for all $\omega' \in \Gamma^*$ if $\langle p, \gamma \rangle \xrightarrow{L} \langle q, \omega \rangle$ and $\omega' \in L$. Given a set of configurations C , we define $cpre^*(C) = \{c' \mid \exists c \in C : c' \Rightarrow_c^* c\}$ and $cpost^*(C) = \{c' \mid \exists c \in C : c \Rightarrow_c^* c'\}$ to be the (possibly infinite) set of pre- and post- images of C , respectively.

A CPDS extends a PDS by further associating each transition with regular languages over the stack that specify conditions under which the transition can be applied

Definition 2.7 A conditional weighted pushdown system \mathcal{W}_c is a triplet (\mathcal{P}_c, S, f_c) , where $\mathcal{P}_c = (Q, \Gamma, C, \Delta_c, q_0, \omega_0)$ is a CPDS, $S = (D, \oplus, \otimes, \bar{0}, \bar{1})$ is a bounded idempotent semiring, and $f_c : \Delta \rightarrow D$ is a weight assignment function.

Li et al. lifted the MOVP problem in Definition 2.4 to CWPDSs by replacing the underlying PDSs with CPDSs, denoted by $MOVP_c(S, T)$, and showed that model checking problems can be reduced to that on WPDSs.

Given a regular language L , a condition automaton $A = (S, \Gamma, \delta, \hat{s}, F)$ with respect to L is a total deterministic finite automaton that recognizes L , where S is the set of states, Γ is the input alphabet, $\delta : S \times \Gamma \rightarrow S$ is the transition function, \hat{s} is the initial state, and $F \subseteq S$ is the set of final states. Let $C = \{L_1, \dots, L_n\}$ in \mathcal{P}_c , and let L^R denote the set of reversed words of words from L .

$\prod_{1 \leq i \leq n} A_i$ is defined as a Cartesian product of all condition automata A_i with respect to conditions L_i^R , and it is not hard to see being total and deterministic, i.e., a condition automaton. We denote by s_i the i^{th} component of a state \mathbf{s} in \hat{S} .

Let $A_i = (S_i, \Gamma, \delta_i, \hat{s}_i, F_i)$ for $1 \leq i \leq n$, and let $\prod_{1 \leq i \leq n} A_i = (\hat{S}, \Gamma, \hat{\delta}, \hat{s}, \hat{F})$, where $\hat{\delta}(\mathbf{s}, \gamma) = (\delta_1(s_1, \gamma), \dots, \delta_n(s_n, \gamma))$ for any $\mathbf{s} \in$

\hat{S} , $\hat{s} = (\hat{s}_1, \dots, \hat{s}_n)$, and $\hat{F} = \{(s_1, \dots, s_n) \mid \exists 1 \leq i \leq n : s_i \in F_i\}$. The authors proposed an *offline algorithm* TRANS that translates a CWPDS \mathcal{W}_c to a WPDS \mathcal{W}_t , by synchronizing the underlying pushdown system and $\prod_{1 \leq i \leq n} A_i$. TRANS translates transitions δ in the left to δ' in the right

$$\begin{aligned} \langle p, \gamma \rangle \xrightarrow{A_i} \langle q, \epsilon \rangle & \quad \langle p, (\gamma, \mathbf{r}) \rangle \leftrightarrow \langle q, \epsilon \rangle \\ \langle p, \gamma \rangle \xrightarrow{A_i} \langle q, \gamma' \rangle & \quad \Rightarrow \quad \langle p, (\gamma, \mathbf{r}) \rangle \leftrightarrow \langle q, (\gamma', \mathbf{r}) \rangle \\ \langle p, \gamma \rangle \xrightarrow{A_i} \langle q, \gamma' \gamma'' \rangle & \quad \langle p, (\gamma, \mathbf{r}) \rangle \leftrightarrow \langle q, (\gamma', \mathbf{t})(\gamma'', \mathbf{r}) \rangle \end{aligned}$$

where $\mathbf{r} \in \hat{S}$, $\mathbf{r}_i \in F_i$ and $\hat{\delta}(\mathbf{r}, \gamma'') = \mathbf{t}$, and $f(\delta') = f_c(\delta)$.

A configuration $\langle p, (\gamma_n, \mathbf{r}_n)(\gamma_{n-1}, \mathbf{r}_{n-1}) \dots (\gamma_0, \mathbf{r}_0) \rangle$ of \mathcal{W}_t is consistent if $\mathbf{r}_0 = \hat{\mathbf{r}}$ and $\hat{\delta}(\mathbf{r}_i, \gamma_i) = \mathbf{r}_{i+1}$ for each $1 \leq i \leq n-1$. Let $\rho(\langle p, (\gamma_n, \mathbf{r}_n)(\gamma_{n-1}, \mathbf{r}_{n-1}) \dots (\gamma_0, \mathbf{r}_0) \rangle) = \langle p, \gamma_n \dots \gamma_0 \rangle$, and let ρ^{-1} be the inverse. The computation of \mathcal{W}_t over consistent configurations and the computation of \mathcal{W}_c are bisimilar, and as such the model checking problems on \mathcal{W}_c is reduced to \mathcal{W}_t over consistent configurations. Note that the time complexity increases from the unconditional case by a factor that is polynomial in the size of \hat{S} , and $|\hat{S}|$ can be exponentially large.

3. An On-The-Fly Algorithm for CWPDSs

The offline algorithm synchronizes the underlying pushdown system and condition automata, and generates the entire space of transitions that may explode due to the potential blow-up of $|\hat{S}|$. However, a large part of the transitions in \mathcal{W}_t may never be used in answering the model checking enquiry. We present an on-the-fly algorithm to compute $MOVP_c(C, T)$ that synchronizes the two computing machineries on-demand while computing post-images. An on-the-fly algorithm in terms of computing pre-images can be similarly constructed. Given a \mathcal{P} -automaton \mathcal{A}_C that recognizes a regular configuration C .

The on-the-fly model checking algorithm consists in consecutive four steps,

- (1) building a \mathcal{P} -automaton $\mathcal{A}_{C'}$ that recognizes consistent configurations C' of \mathcal{W}_t such that $C' = \{\rho^{-1}(c) \mid c \in C\}$, as given in Algorithm 1, and
- (2) applying the saturation procedure to computing a \mathcal{P} -automaton $\mathcal{A}_{cpost^*(C')}$ that accepts $cpost^*(C')$, by synchronizing the pushdown system and condition automata on-demand, as given in Algorithm 2, and
- (3) building a \mathcal{P} -automaton $\mathcal{A}_{T'}$ that recognizes consistent configurations T' of \mathcal{W}_t such that $T' = \{\rho^{-1}(c) \mid c \in T\}$, as done in Step (1), and computing a \mathcal{P} -automaton \mathcal{B} by $\mathcal{B} = \mathcal{A}_{T'} \cap \mathcal{A}_{cpost^*(C')}$, and
- (4) computing $MOVP_c(C, T)$ by running Algorithm 3 on \mathcal{B} , which is based on Algorithm 4 in Fig. 19 of Ref. [4].

Note that, \mathcal{A}_C and $\mathcal{A}_{T'}$ are unweighed automata, and $\mathcal{A}_{cpost^*(C')}$ is a weighted automaton. They all accept consistent configurations of \mathcal{W}_i , although \mathcal{W}_i is not explicitly constructed at all in the on-the-fly algorithm. The intersection of $\mathcal{A}_{T'}$ and $\mathcal{A}_{cpost^*(C')}$ is done the same way as standard, except that the weights of $\mathcal{A}_{cpost^*(C')}$ are carried over to \mathcal{B} .

In the sequel, we fix a CWPDS $\mathcal{W}_c = (\mathcal{P}_c, \mathcal{S}, f)$, where $\mathcal{P}_c = (Q, \Gamma, C, \Delta_c, q_0, \gamma_0)$, $C = \{L_1, \dots, L_n\}$ and $\mathcal{S} = (D, \oplus, \otimes, \bar{0}, \bar{1})$. Without loss of generality, we assume $\gamma_0 \in \Gamma$. Let $A = \{A_1, \dots, A_n\}$ where $A_i = (S_i, \Gamma, \delta_i, \hat{s}_i, F_i)$ is the condition automaton with respect to L_i^R for $1 \leq i \leq n$.

Algorithm 1 *BuildInitial*(\mathcal{A}_C)

Require: a \mathcal{P} -automaton $\mathcal{A}_C = (Q', \Gamma, \rightarrow, Q, F)$ that recognizes a regular set of configuration $C \subseteq Q \times \Gamma^*$, and \mathcal{A}_C has no transitions into Q states and has no ε -transitions.

Ensure: a \mathcal{P} -automaton $\mathcal{A}_{C'} = (H, \Gamma', \rightarrow', Q, F')$ that recognizes C' where $C' = \{\rho^{-1}(c) \mid c \in C\}$.

```

1:  $F' := \emptyset; \rightarrow' := \emptyset; visited := \emptyset$ 
2:  $\hat{s} := (\hat{s}_1, \hat{s}_2, \dots, \hat{s}_n)$ 
3: for each  $q \in F$  do
4:    $F' := F' \cup \{(q, \hat{s})\}$ 
5:  $H := F'$ 
6:  $workset := F'$ 
7: while  $workset \neq \emptyset$  do
8:   select and remove a state  $u = (p, s)$  from  $workset$ 
9:    $visited := visited \cup \{(p, s)\}$ 
10:  for each  $(q, \gamma, p) \in \rightarrow$  do
11:     $t := \hat{\delta}(s, \gamma)$ 
12:     $H := H \cup \{(q, t)\}$ 
13:     $\rightarrow' := \rightarrow' \cup \{(q, t), (\gamma, s), (p, s)\}$ 
14:    if  $(p, t) \notin visited$  then
15:       $workset := workset \cup \{(q, t)\}$ 
16:  for each  $((q, t), (\gamma, s), (p, s)) \in \rightarrow'$  do
17:    if  $q \in Q$  then
18:       $\rightarrow' := \rightarrow' \cup \{(q, (\gamma, s), (p, s))\}$ 
19: return  $(H, \Gamma', \rightarrow', Q, F')$ 
    
```

Algorithm 1 gives a procedure that, given \mathcal{A}_C that accepts a regular set C of configurations, builds the \mathcal{P} -automaton $\mathcal{A}_{C'}$ such that $C' = \{\rho^{-1}(c) \mid c \in C\}$. Since \mathcal{A}_C has no ε -transitions, it is not hard to conclude by the algorithm construction that $\mathcal{A}_{C'}$ has no transitions into Q states and has no ε -transitions. By properly choosing data structures, all the needed membership test, removal operations, etc., take constant time. Let L be the longest length of descending chains and t_{op} be the time of performing binary operators of the semiring \mathcal{S} . By a simple analysis over line 6 to 12, Algorithm 1 takes $O(|\rightarrow| \cdot |Q| \cdot |\hat{S}| \cdot |t_{op}| \cdot L)$ time and space.

Algorithm 2 gives the on-the-fly procedure for computing $cpost^*(C')$. The synchronization points for condition checking with the original procedure for computing $post^*$ are highlighted by underlines. Note that the procedure for computing and propagating weights is kept unchanged. It is not hard to conclude by the algorithm construction that $\mathcal{A}_{cpost^*(C')}$ has no transitions into Q states, given the above invariant on Algorithm 1. The algorithm extends the one for efficiently computing A_{post^*}

Algorithm 2 Computing $\mathcal{A}_{cpost^*(C')}$ On The Fly

Require: a \mathcal{P} -automaton $\mathcal{A}_C = (Q', \Gamma, \rightarrow, Q, F)$ that recognizes a regular set of configuration $C \subseteq Q \times \Gamma^*$, and \mathcal{A}_C has no transitions into Q states and has no ε -transitions.

Ensure: a \mathcal{P} -automaton $\mathcal{A}_{cpost^*(C')} = (H', \Gamma', \rightarrow', Q, F')$ that recognizes $cpost^*(C)$ of \mathcal{W}_i , where $C' = \{\rho^{-1}(c) \mid c \in C\}$; and a mapping $l : (\rightarrow') \rightarrow D$.

```

1: procedure  $UPDATE(t, v)$ 
2:   begin
3:      $newValue := l(t) \oplus v$ 
4:      $changed := (newValue \neq l(t))$ 
5:     if  $changed$  then
6:        $workset := workset \cup \{t\}$ 
7:        $l(t) := newValue$ 
8:   end
9:  $\mathcal{A}_{C'} := BuildInitial(\mathcal{A}_C)$ 
10: Let  $\mathcal{A}_{C'} = (H, \Gamma', \rightarrow, Q, F')$ 
11:  $\rightarrow' := \rightarrow_0; workset := \rightarrow_0; l := \lambda t. \bar{0}$ 
12: for all  $t \in \rightarrow_0$  do  $l(t) := \bar{1}$ 
13:  $H' := H$ 
14: while  $workset \neq \emptyset$  do
15:   take and remove a transition  $t = (p, (\gamma, \mathbf{u}), q)$  from  $workset$ 
16:   if  $\gamma \neq \varepsilon$  then
17:     for all  $r = \langle p, \gamma \rangle \xrightarrow{A_i} \langle p', \varepsilon \rangle \in \Delta_c$  with  $A_i = (S_i, \Gamma, \delta_i, \hat{s}_i, F_i)$  do
18:       if  $\mathbf{u}_i \in F_i$  then
19:          $UPDATE(\langle p', \varepsilon, q \rangle, l(t) \otimes f(r))$ 
20:       for all  $r = \langle p, \gamma \rangle \xrightarrow{A_i} \langle p', \gamma' \rangle \in \Delta_c$  with  $A_i = (S_i, \Gamma, \delta_i, \hat{s}_i, F_i)$  do
21:         if  $\mathbf{u}_i \in F_i$  then
22:            $UPDATE(\langle p', (\gamma', \mathbf{u}), q \rangle, l(t) \otimes f(r))$ 
23:         for all  $r = \langle p, \gamma \rangle \xrightarrow{A_i} \langle p', \gamma' \gamma'' \rangle \in \Delta_c$  with  $A_i = (S_i, \Gamma, \delta_i, \hat{s}_i, F_i)$  do
24:           if  $\mathbf{u}_i \in F_i$  then
25:              $let \mathbf{z} = \hat{\delta}(\mathbf{u}, \gamma'')$ 
26:              $H' := H' \cup \{h_{(p', \gamma', \mathbf{z})}\}$ 
27:              $UPDATE(\langle p', (\gamma', \mathbf{z}), (h_{(p', \gamma', \mathbf{z})}, \mathbf{z}), \bar{1} \rangle)$ 
28:              $UPDATE(\langle h_{(p', \gamma', \mathbf{z})}, \mathbf{z} \rangle, (\gamma'', \mathbf{u}), q, l(t) \otimes f(r))$ 
29:             if  $changed$  then
30:               for all  $t' = \langle p'', \varepsilon, (h_{(p', \gamma', \mathbf{z})}, \mathbf{z}) \rangle$  do
31:                  $UPDATE(\langle p'', (\gamma'', \mathbf{u}), q \rangle, l(t) \otimes f(r) \otimes l(t'))$ 
32:           else
33:             for all  $t' = \langle q, (\gamma', \mathbf{w}), q' \rangle \in \rightarrow$  do
34:                $UPDATE(\langle p, (\gamma', \mathbf{w}), q' \rangle, l(t') \otimes l(t))$ 
35:   return  $((H', \Gamma', \rightarrow', Q, F'), l)$ 
    
```

(Algorithm 2 of Section 3) in Ref. [5], and the complexity increased by a factor of $n_1 = |\hat{S}| \cdot |t_{op}| \cdot L$. Algorithm 2 takes $O((|Q| \cdot |\Delta| \cdot (|H \setminus Q| + n_0) + |Q| \cdot |\rightarrow_0|) \cdot n_1)$ time and space, where n_0 is the number of different pairs (p, γ) such that there exists a rule $\langle p, \gamma \rangle \xrightarrow{A_i} \langle p', \gamma' \gamma'' \rangle$ in Δ_c .

Remark 3.1 By the construction of Algorithm 1, we have that $\mathbf{s} = \mathbf{s}'$ for each transition $(\rightarrow, (\gamma, \mathbf{s}'), (p, \mathbf{s}))$ of $\mathcal{A}_{C'}$. Based on this invariant on $\mathcal{A}_{C'}$ and the construction of Algorithm 2, the same fact applies to transitions of $\mathcal{A}_{cpost^*(C')}$. It indicates that we can further enhance the performance of the algorithm by first removing the second components of Γ' in Algorithm 1 and 2, and then directly compute $\mathcal{B} = \mathcal{A}_T \cap \mathcal{A}_{cpost^*(C')}$, instead of computing

Algorithm 3 *GenValue*(\mathcal{B})

Require: a \mathcal{P} -automaton $\mathcal{B} = (H, \Gamma, \rightarrow, Q, F)$ that has no transitions into Q states, and a mapping l from \rightarrow to D .

Ensure: a value $d \in D$

```

1: for each  $u \in H \setminus (Q \cup F)$  do
2:    $V_u := \bar{0}$ 
3: for each  $u \in F$  do
4:    $V_u := \bar{1}$ 
5:  $workset := F$ 
6: while  $workset \neq \emptyset$  do
7:   select and remove a state  $u$  from  $workset$ 
8:   for each  $t = (s, \gamma, u) \in \rightarrow$  do
9:      $newValue := V_s \oplus (V_u \otimes l(t))$ 
10:    if  $newValue \neq V_s$  then
11:       $V_s := newValue$ 
12:       $workset := workset \cup \{s\}$ 
13:  $val := \bar{0}$ 
14: for each  $u \in Q$  do
15:    $val := val \oplus V_u$ 
return  $val$ 
    
```

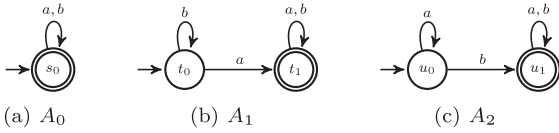


Fig. 2 Condition automata A_0, A_1 and A_2 .

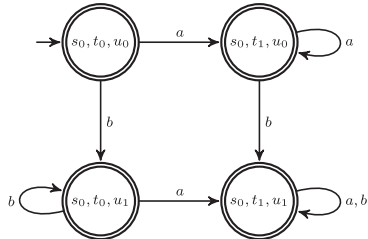


Fig. 3 $\prod_{1 \leq i \leq n} A_i$.

$\mathcal{A}_{T'}$ and $\mathcal{B} = \mathcal{A}_{T'} \cap \mathcal{A}_{cpost^*(C')}$. Note that this optimisation would be crucial to practical efficiency. Without this optimisation, $\mathcal{A}_{T'}$ has to be computed, which, however, could easily explode.

Example 3.2 Consider a CPDS $\mathcal{P}_c = (Q, \Gamma, \{A_0, A_1, A_2\}, \Delta_c, p_0, a)$, where $Q = \{p_0, p_1, p_2\}$, $\Gamma = \{a, b\}$, and the reverse of regular languages that represent regular conditions are directly described by condition automata A_0, A_1, A_2 , as shown in **Fig. 2**, where $L(A_0) = \Gamma^*$, $L(A_1) = b^*a\Gamma^*$ and $L(A_2) = a^*b\Gamma^*$, and their Cartesian product A is depicted in **Fig. 3**.

Transitions are as follows

$$\Delta_c = \left\{ \begin{array}{l} \langle p_0, a \rangle \xrightarrow{A_0} \langle p_0, ba \rangle \\ \langle p_0, b \rangle \xrightarrow{A_0} \langle p_0, bb \rangle \\ \langle p_0, b \rangle \xrightarrow{A_0} \langle p_0, \varepsilon \rangle \\ \langle p_0, b \rangle \xrightarrow{A_1} \langle p_1, b \rangle \\ \langle p_0, a \rangle \xrightarrow{A_2} \langle p_2, a \rangle \end{array} \right\}$$

Given the initial \mathcal{P} -automaton \mathcal{A}_S in **Fig. 4**, we are interested in knowing whether states p_1 and p_2 are reachable or not, i.e., there exists $\omega \in \Gamma^*$ such that $p_i \xrightarrow{\omega}^* q$ for some final state q in the output \mathcal{P} -automaton, where $i \in \{0, 1\}$.

We consider a CWPDS with \mathcal{P}_c as the underlying CPDS and

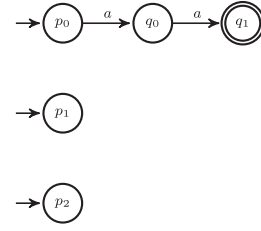


Fig. 4 The Initial \mathcal{P} -automaton \mathcal{A}_S .

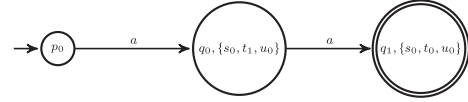


Fig. 5 $\mathcal{A}_{S'}$ after applying Algorithm 1.

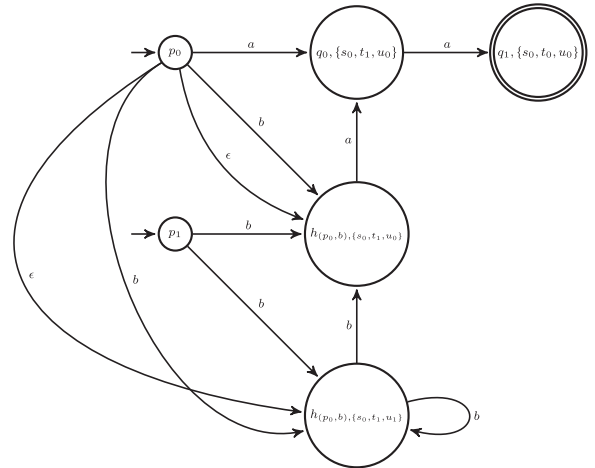


Fig. 6 $\mathcal{A}_{cpost^*(S')}$ after applying Algorithm 2.

the semiring $(\{\bar{0}, \bar{1}\}, \oplus, \otimes, \bar{0}, \bar{1})$. The result after applying Algorithm 1 and Algorithm 2 are shown in **Fig. 5** and **Fig. 6**, respectively. Given the initial \mathcal{P} -automaton, \mathcal{P}_c generates configurations $\langle p_0, b^*aa \rangle$. We know p_1 is reachable and p_2 is not. The result is correctly shown by the output automaton.

4. Experiments

4.1 Experimental Configuration

We implement all the algorithms in Java (Java Development Kit 1.7). Our package extends the weighted pushdown system library of jMoped [6]. We use the Java automaton library developed by Anders Moller [3] to represent regular conditions associated with transitions by finite automata. The library provides facilities of a fast manipulation and easy translation from regular expressions to deterministic automata. Experiments are conducted on a system with Intel®Xeon®CPU E5-2690 0 @2.9 GHz, Windows Server 2008 R2 Standard with 32 GB RAM.

Comparison between the offline algorithm and the on-the-fly algorithm is conducted with a big realistic example obtained from a study of reachability analysis of HTML5 parser specifications [2], in which specifications are translated to CPDSs. Some details of the input CPDS are:

- Initial configuration $S = \{\langle 323, X \rangle\}$.
- Size of stack alphabet $|\Gamma| = 25$.
- Number of control locations $|Q| = 487$.

Table 1 A Comparison of offline and on-the-fly algorithms.

C	$\mathcal{A}_{\text{post}^*}$	Offline algorithm		On-the-fly algorithm	
		Time	Memory	Time	Memory
2	31 414	1 s	61 268	0.798 s	56 949
5	31 412	8 s	1 147 307	0.82 s	57 266
10	109 960	∞	∞	1 s	96 678
15	211 368	∞	∞	2 s	43 877
20	418 363	∞	∞	2 s	138 844
60	3 445 895	∞	∞	12 s	760 871
70	3 244 153	∞	∞	12 s	877 967
80	3 331 394	∞	∞	15 s	715 939
118	32 075 528	∞	∞	3 m 46 s	7 410 874

- Number of pushdown transitions $|\Delta| = 19679$.
- Each condition L_i in the system is described by a regular expression, e.g.,
 $((\text{Div} | (\text{Optgroup} | (\text{Option} | (\text{P} | (\text{Rp} | \text{Ruby})))))) * \text{Li} (@)$,
 where @ means any word. The corresponding condition automaton A_i may contain up to 6 states.

4.2 Experimental Results

Our experimental results are given in **Table 1**, which compares the space and time efficiency of the offline and on-the-fly algorithms. The first column $|C|$ shows the number of conditions in CPDSs, for which we first remove conditions of all transitions and then increase conditions randomly. Recall that the number of conditions, i.e., the dimension of the product automaton, is the factor causing the state blow-up. The second column gives the size of the output \mathcal{P} -automaton, in terms of the number of transitions, after applying Algorithm 2. The runtime and memory consumption are estimated in seconds (s) and Kilobyte (KB), given in the next columns, and ∞ means either taking longer than 24 hours or getting out-of-memory.

When the number of conditions $|C|$ is up to 5, the offline algorithm is still able to work, but the resource and overhead cost taken is at least two times more than that of the on-the-fly algorithm. When the number of conditions $|C|$ is over 5, theoretically the product of condition DFA could explode exponentially as much as $6^{|C|}$ in size, where 6 indicates the maximum number of states in a single condition DFA. It is not surprising that the offline algorithm run out of memory. Table 1 proves efficiency of the on-the-fly algorithm towards performance. The last row of Table 1 corresponds to the runtime profile of the entire example. It shows that analysing the big example up to step (2) takes about 4 minutes and 7 GB memory. Note that, the memory usage takes into account garbage collection. We allocated 24 GB to the Java virtual machine for running the analysis.

4.3 Discussion

Taking the semiring containing only unit elements $\bar{0}$ and $\bar{1}$, we computed $\text{MOVPC}_c(S, T)$ for target configurations T generated from the study of reachability analysis of HTML5 parser specifications [2]. Since the model checking problem is EXPTIME-complete, the runtime performance of the analysis at large depends on practical instances. For instance, it always took our package less than 2 secs to analyse targets in the form of $\{\langle p, \gamma \omega \rangle \mid \omega \in \Gamma^*, \gamma \in \text{Sym}\}$, denoted by $(p, \text{Sym}\Gamma^*)$, where $p \in Q$ is some control location, and $\text{Sym} \subseteq \Gamma$ is a set of stack symbols. For

those form of targets, we can safely halt the analysis once any head of the target configurations is reachable during the saturation. However, our analysis run out of memory for other form of targets during the step (3), i.e, making intersection of $\mathcal{A}_{\text{post}^*} \cap A_T$ (which may result in a quadratic increase in the automaton size), for which the garbage collector keeps up hard work to gain more space.

The algorithm in Ref. [2] computes backward saturation of target configurations and our analysis computes forward saturation from the source configurations. The advantage of our analysis is that $\mathcal{A}_{\text{post}^*}(C)$ is computed once and then reused for analysing all target configurations. In contrast, the analysis in Ref. [2] need to compute $\mathcal{A}_{\text{pre}^*}(T)$ for each target T . However, backward saturation has an inherit merit of possibly exploring smaller state space as opposed to forward saturation. We expect that our analysis would likely excel at analysing problem instances generated by the analysis and verification of object-oriented programs as originally discussed in Ref. [1], which do not contain as many as distinguished regular conditions as those in Ref. [2], and targets of interest in program analysis and verification are in those simple forms as discussed above.

5. Conclusion

We present an on-the-fly algorithm for model checking CW-PDSs, to tackle the space explosion problem caused by its original offline algorithm. The idea is to interleave condition checking with the saturation procedures for computing post-images of regular configurations. We implemented the on-the-fly algorithm in terms of computing post-images and conducted experiments on realistic big examples from the study of reachability analysis of HTML5 parser specification and compatibility checking [2]. Our experimental results show that the on-the-fly algorithm drastically outperforms the offline algorithm regarding both space and time efficiency in practice. It would be interesting to apply our model checker to analysis and verification problems of object-oriented programs like Java.

Acknowledgments We would like to thank Yasuhiko Minamide for sharing us with their model and test data generated from HTML5 parser specifications. We thank Mizuhito Ogawa for comments and support. This work is partially supported by Kakenhi 25730039.

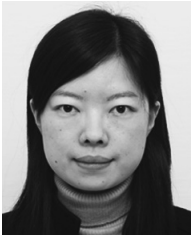
References

- [1] Li, X. and Ogawa, M.: Conditional weighted pushdown systems and applications, *Proc. ACM SIGPLAN 2010 Workshop on Partial evaluation and Program Manipulation — PEPM '10*, p.141 (online), DOI: 10.1145/1706356.1706382 (2010).
- [2] Minamide, Y. and Mori, S.: Reachability Analysis of the HTML5 Parser Specification and its Application to Compatibility Testing, *18th International Symposium on Formal Methods* (2012).
- [3] Moller, A.: dk.brics.automaton — Finite-State Automata and Regular Expressions for Java (2010).
- [4] Reps, T., Schwoon, S., Jha, S. and Melski, D.: Weighted pushdown systems and their application to interprocedural dataflow analysis, *Science of Computer Programming*, Vol.58, No.1-2, pp.206–263 (online), DOI: 10.1016/j.scico.2005.02.009 (2005).
- [5] Schwoon, S.: Model-Checking Pushdown Systems, Ph.D. Thesis, Technische Universität München (2002).
- [6] Suwimonteerabuth, D.: Reachability in Pushdown Systems: Algorithms and Applications, PhD Thesis, Technische Universität München (2009).

- [7] Xin Li and Hua, V.L.T.: Generating Stack-based Access Control Policies, *arXiv.org*, pp.1-18 (online), available from <http://arxiv.org/abs/1307.2964v2> (2013).



Hua Vy Le Thanh received her M.S. in 2012 from Japan Advanced Institute of Science and Technology. Her research interests include embedded software, security and program verification methodology, such as model checking and program analysis.



Xin Li received her M.S. from Shanghai Jiao Tong University in 2004 and her Ph.D from JAIST (Japan Advanced Institute of Science and Technology) in 2007, respectively. She worked at JAIST from 2007 to 2013, and is presently with the University of Tokyo. Her current research interests are program analysis and verification.