

Pushdown Model Checking for Malware Detection[★]

Fu Song and Tayssir Touili

LIAFA, CNRS and Univ. Paris Diderot, France.

E-mail: {song,touili}@liafa.jussieu.fr

Abstract. The number of malware is growing extraordinarily fast. Therefore, it is important to have efficient malware detectors. Malware writers try to obfuscate their code by different techniques. Many of these well-known obfuscation techniques rely on operations on the stack such as inserting dead code by adding useless push and pop instructions, or hiding calls to the operating system, etc. Thus, it is important for malware detectors to be able to deal with the program's stack. In this paper we propose a new model-checking approach for malware detection that takes into account the behavior of the stack. Our approach consists in : (1) Modeling the program using a Pushdown System (PDS). (2) Introducing a new logic, called SCTPL, to represent the malicious behavior. SCTPL can be seen as an extension of the branching-time temporal logic CTL with variables, quantifiers, and predicates over the stack. (3) Reducing the malware detection problem to the model-checking problem of PDSs against SCTPL formulas. We show how our new logic can be used to precisely express malicious behaviors that could not be specified by existing specification formalisms. We then consider the model-checking problem of PDSs against SCTPL specifications. We reduce this problem to emptiness checking in Symbolic Alternating Büchi Pushdown Systems, and we provide an algorithm to solve this problem. We implemented our techniques in a tool, and we applied it to detect several viruses. Our results are encouraging.

1 Introduction

To identify viruses, existing antivirus systems use either code emulation or signature (pattern) detection. These techniques have some limitations. Indeed, emulation based techniques can only check the program's behavior in a limited time interval, whereas signature based systems are easy to get around. To sidestep these limitations, instead of executing the program or making a syntactic check over it, virus detectors need to use analysis techniques that check the *behavior* (not the syntax) of the program in a *static* way, i.e. without executing it. Towards this aim, we propose in this paper to use *model-checking* for virus detection. Model-checking has already been used for virus detection in [6, 20, 9, 11, 16, 15, 17]. However, these works model the program as a finite state graph (automaton). Thus, they are not able to model the stack of the programs, and cannot track the effects of the push, pop and call instructions. However, as described in [19], many obfuscation techniques rely on operations over the stack. Indeed, many antivirus systems determine whether a program is malicious by checking the calls it

[★] Work partially funded by ANR grant ANR-08-SEGI-006.

makes to the operating system. Hence, several virus writers try to hide these calls by replacing them by push and return instructions [19]. Therefore, it is important to have analysis techniques that can deal with the program stack.

We propose in this paper a novel model-checking technique for malware detection that takes into account the behavior of the stack. Our approach consists in modeling the program using a *pushdown system* (PDS), and defining a new logic, called *SCTPL*, to express the malicious behavior.

Using pushdown systems as program model allows to consider the program stack. In our modeling, the PDS control locations correspond to the program's control points, and the PDS stack mimics the program's execution stack. This allows the PDS to mimic the behavior of the program. This is different from standard program translations to PDSs where the control points of the program are stored in the stack [13, 5]. These standard translations assume that the program follows a standard compilation model, where the return addresses are never modified. We do not make such assumptions since behaviors where the return addresses are modified can occur in malicious code. We only make the assumption that pushes and pops can be done only using *push*, *pop*, *call*, and *return* operations, not by manipulating the stack pointer.

The logic SCTPL that we introduce is an extension of the CTPL logic that allows to use predicates over the stack. CTPL was introduced in [16, 15, 17]. It can be seen as an extension of CTL with variables and quantifiers. In CTPL, propositions can be predicates of the form $p(x_1, \dots, x_n)$, where the x_i 's are free variables or constants. Free variables can get their values from a finite domain. Variables can be universally or existentially quantified. CTPL is as expressive as CTL, but it allows a more succinct specification of the malicious behavior. For example, consider the statement "The value *data* is assigned to some register, and later, the content of this register is pushed onto the stack." This statement can be expressed in CTL as a large formula enumerating all the possible registers:

$$\begin{aligned} & EF (mov(eax, data) \wedge AF push(eax)) \vee \\ & EF (mov(ebx, data) \wedge AF push(ebx)) \vee \\ & EF (mov(ecx, data) \wedge AF push(ecx)) \vee \dots \end{aligned}$$

where every instruction is regarded as a predicate, i.e., $mov(eax, data)$ is a predicate. However, the CTL formula is large for such a simple statement. Using CTPL, this can be expressed by the CTPL formula $\exists r EF (mov(r, data) \wedge AF push(r))$ which expresses in a succinct way that there exists a *register* r such that the above holds. [16, 15, 17] show how this logic is adequate to specify some malicious behaviors. However, CTPL does not allow to specify properties about the stack (which is important for malicious code detection as explained above).

For example, consider Figure 1(a). It corresponds to a critical fragment of the

Email-worm Avron [14] that shows the typical behavior of an email worm: it calls an API function *GetModuleHandleA* with 0 as its parameter. This allows to get the entry address of its own executable so that later, it can infect other files by copying this executable into them. (Parameters to a function in assembly are passed by pushing them onto the stack before a call to the function is made. The code in the called function later retrieves these parameters from the stack.) Using CTPL, we can specify this malicious behavior by the following formula:

$$\exists r_1 \mathbf{EF} \left(\text{mov}(r_1, 0) \wedge \mathbf{EX E} [\neg \exists r_2 \text{mov}(r_1, r_2) \mathbf{U} (\text{push}(r_1) \wedge \mathbf{EX E} [\neg \exists r_3 (\text{push}(r_3) \vee \text{pop}(r_3)) \mathbf{U} \text{call}(\text{GetModuleHandleA})])]) \right) \quad (1)$$

This formula states that there exists a register r_1 assigned by 0 such that the value of r_1 is not modified until it is pushed onto the stack. Later the stack is not changed until function *GetModuleHandleA* is called. This specification can detect the fragment in Figure 1(a). However, a worm writer can easily use some obfuscation techniques in order to escape this specification. For example, let us introduce one push followed by one pop after *push eax* at line l_2 as done in Figure 1(b). By doing so, this fragment keeps the same malicious behavior than the fragment in Figure 1(a). However, it cannot be detected by the above CTPL formula. Since the number of pushes and pops that can be added by the worm writer can be arbitrarily large, it is always possible for worm developers to change their code in order to escape a given CTPL formula.

To overcome this problem, we introduce the SCTPL logic which extends CTPL by predicates over the stack. Such predicates are given by regular expressions over the stack alphabet and some free variables (which can also be existentially and universally quantified). Using our new logic SCTPL, the malicious behavior of Figures 1(a) and (b) can be specified as follows:

$$\psi = \exists r_1 \mathbf{EF} \left(\text{mov}(r_1, 0) \wedge \mathbf{EX E} [\neg \exists r_2 \text{mov}(r_1, r_2) \mathbf{U} (\text{push}(r_1) \wedge \mathbf{EX E} [\neg (\text{push}(r_1) \vee (\exists r_3 (\text{pop}(r_3) \wedge r_1 I^*))) \mathbf{U} (\text{call}(\text{GetModuleHandleA}) \wedge r_1 I^*)])]) \right) \quad (2)$$

where $r_1 I^*$ is a regular predicate expressing that the topmost symbol of the stack is r_1 . The SCTPL formula ψ states that there exists a register r_1 assigned by 0 such that the value of r_1 is not changed until it is pushed onto the stack. Then, r_1 is never pushed onto the stack again nor popped from it until the function *GetModuleHandleA* is called. When this call is made, the topmost symbol of the stack has to be r_1 . This ensures that *GetModuleHandleA* is called with 0 as parameter. This specification can detect both fragments in Figure 1, because it allows to specify the content of the stack when *GetModuleHandleA* is called. Note that it is important to use pushdown systems as model in order to have specifications with predicates over the stack.

The main contributions of this paper are:

1. We present a new technique to translate a binary program into a pushdown system that mimics the program's behavior (a malicious program is usually an executable, i.e., a binary program). Our translation is different from standard program translations to PDSs that need to assume that the program follows a standard compilation

```
l1: mov eax,0
l2: push eax
l3: call ds:GetModuleHandleA
```

(a)

```
l'1: mov eax,0
l'2: push eax
l'3: push ebx
l'4: pop ebx
l'5: call ds:GetModuleHandleA
```

(b)

Fig. 1. (a) Worm fragment; (b) Obfuscated fragment.

model, where the return addresses are never modified. Our translation does not need to make this assumption since malicious code may have a non standard form.

2. We introduce the SCTPL logic and show how it can be used to efficiently and precisely characterize malicious behaviors.
3. We propose an algorithm for model checking pushdown systems against SCTPL specifications. We reduce this problem to checking emptiness in Symbolic Alternating Büchi Pushdown Systems (SABPDS) and we propose an algorithm to solve this emptiness problem.
4. We implemented our techniques in a tool that we successfully applied to detect several viruses.

Related work. Model-checking and static analysis techniques have been applied to detect malicious behaviors e.g. in [6, 20, 9, 11, 16, 15, 17]. However, all these works are based on modeling the program as a finite-state system, and thus, they miss the behavior of the stack. As we have seen, being able to track the stack is important for many malicious behaviors. [7] use tree automata to represent a set of malicious behaviors. However, [7] cannot specify predicates over the stack content.

[19] keeps track of the stack by computing an abstract stack graph which finitely represents the infinite set of all the possible stacks for every control point of the program. Their technique can detect only obfuscated calls and obfuscated returns. Using SCTPL, we are able to detect more malicious behaviors.

[18] performs context-sensitive analysis of *call* and *ret* obfuscated binaries. They use abstract interpretation to compute an abstraction of the stack. We believe that our techniques are more precise since we do not abstract the stack. Moreover, the techniques of [18] were only tried on toy examples, they have not been applied for malware detection.

[5] uses pushdown systems for binary code analysis. However, [5] has not been applied for malware detection. Moreover, the translation from programs to PDSs in [5] assumes that the program follows a standard compilation model where calls and returns match. Several malicious behaviors do not follow this model. Our translation from a control flow graph to a PDS does not make this assumption.

[10] defines a language for specifying malicious behavior in terms of dependences between system calls. Compared to SCTPL, the specification language of [10] does not take the stack into account and is only able to express safety properties (no CTL like properties), whereas SCTPL does. On the other hand, [10] is able to automatically derive the malicious specifications by comparing the execution behavior of a known malware against the execution behaviors of a set of benign programs. It would be interesting to see if their techniques can be extended to automatically derive SCTPL specifications of malicious behaviors.

LTL or CTL model-checking with regular predicates over the stack was considered in [12, 21]. These works do not consider variables and quantifiers.

Outline. We give our formal model in Section 2. In Section 3, we introduce our SCTPL logic. Our SCTPL model checking algorithm for pushdown systems is given in Section 4. The experiments we made for malware detection are reported in Section 5.

2 Formal model: Pushdown Systems

We model a binary code by a pushdown system (PDS). In our modeling, the PDS control locations correspond to the program's control points, and the PDS stack mimics the program's execution stack. This is different from standard program translations to PDSs where the control points of the program are stored in the stack [13, 5]. These standard translations assume that the program follows a standard compilation model, where the return addresses are never modified. We do not make such assumptions since behaviors where the return addresses are modified can occur in malicious code. We only make the assumption that pushes and pops can be done only using *push*, *pop*, *call*, and *return* operations, not by manipulating the stack pointer.

Formally, a *Pushdown System* (PDS) is a tuple $\mathcal{P} = (P, \Gamma, \Delta, \#)$, where P is a finite set of control locations, Γ is the stack alphabet, $\Delta \subseteq (P \times \Gamma) \times (P \times \Gamma^*)$ is a finite set of transition rules, and $\# \in \Gamma$ is the bottom stack symbol. A configuration of \mathcal{P} is $\langle p, \omega \rangle$, where $p \in P$ and $\omega \in \Gamma^*$. If $((p, \gamma), (q, \omega)) \in \Delta$, we write $\langle p, \gamma \rangle \hookrightarrow \langle q, \omega \rangle$. For technical reasons, we assume that the bottom stack symbol $\#$ is never popped from the stack, i.e., there is no transition rule of the form $\langle p, \# \rangle \hookrightarrow \langle q, \omega \rangle \in \Delta$.

The successor relation $\rightsquigarrow_{\mathcal{P}} \subseteq (P \times \Gamma^*) \times (P \times \Gamma^*)$ is defined as follows: if $\langle p, \gamma \rangle \hookrightarrow \langle q, \omega \rangle$, then $\langle p, \gamma\omega' \rangle \rightsquigarrow_{\mathcal{P}} \langle q, \omega\omega' \rangle$ for every $\omega' \in \Gamma^*$. For every configuration $c, c' \in P \times \Gamma^*$, c is a successor of c' iff $c \rightsquigarrow_{\mathcal{P}} c'$. A path is a sequence of configurations c_0, c_1, \dots s.t. $c_i \rightsquigarrow_{\mathcal{P}} c_{i+1}$ for every $i \geq 0$.

3 Malicious behavior Specification

In this section, we introduce the Stack Computation Tree Predicate Logic (SCTPL), the formalism we use to specify malicious behavior.

3.1 Environments, predicates and regular expressions

From now on, we fix the following notations. Let $\mathcal{X} = \{x_1, x_2, \dots\}$ be a finite set of variables ranging over a finite domain \mathcal{D} . Let $B : \mathcal{X} \cup \mathcal{D} \rightarrow \mathcal{D}$ be an environment function that assigns a value $c \in \mathcal{D}$ to each variable $x \in \mathcal{X}$, and such that $B(c) = c$ for every $c \in \mathcal{D}$. $B[x \leftarrow c]$ denotes the environment function such that $B[x \leftarrow c](x) = c$ and $B[x \leftarrow c](y) = B(y)$ for every $y \neq x$. $Abs_x(B)$ is the set of all the environments B' s.t. for every $y \neq x$, $B'(y) = B(y)$. Let \mathcal{B} be the set of all the environment functions.

Let $AP = \{a, b, c, \dots\}$ be a finite set of atomic propositions, $AP_{\mathcal{X}}$ be a finite set of atomic predicates of the form $b(\alpha_1, \dots, \alpha_m)$ such that $b \in AP$, $\alpha_i \in \mathcal{X} \cup \mathcal{D}$ for every i , $1 \leq i \leq m$, and $AP_{\mathcal{D}}$ be a finite set of atomic predicates of the form $b(\alpha_1, \dots, \alpha_m)$ such that $b \in AP$ and $\alpha_i \in \mathcal{D}$ for every i , $1 \leq i \leq m$.

Let $\mathcal{P} = (P, \Gamma, \Delta, \#)$ be a PDS s.t. $\Gamma \subseteq \mathcal{D}$. Let \mathcal{R} be a finite set of regular variable expressions e over $\mathcal{X} \cup \Gamma$ defined by:

$$e ::= \emptyset \mid \epsilon \mid a \in \mathcal{X} \cup \Gamma \mid e + e \mid e \cdot e \mid e^*$$

The language $L(e)$ of a regular variable expression e is a subset of $P \times \Gamma^* \times \mathcal{B}$ defined inductively as follows: $L(\emptyset) = \emptyset$; $L(\epsilon) = \{(\langle p, \epsilon \rangle, B) \mid p \in P, B \in \mathcal{B}\}$; $L(x)$, where $x \in \mathcal{X}$

is the set $\{(\langle p, \gamma \rangle, B) \mid p \in P, \gamma \in \Gamma, B \in \mathcal{B} : B(x) = \gamma; L(\gamma), \text{ where } \gamma \in \Gamma \text{ is the set } \{(\langle p, \gamma \rangle, B) \mid p \in P, B \in \mathcal{B}\}; L(e_1 + e_2) = L(e_1) \cup L(e_2); L(e_1 \cdot e_2) = \{(\langle p, \omega_1 \omega_2 \rangle, B) \mid (\langle p, \omega_1 \rangle, B) \in L(e_1); (\langle p, \omega_2 \rangle, B) \in L(e_2)\}; \text{ and } L(e^*) = \{(\langle p, \omega^* \rangle, B) \mid (\langle p, \omega \rangle, B) \in L(e)\}$. E.g., $(\langle p, \gamma_1 \gamma_1 \gamma_2 \rangle, B)$ is an element of $L(x^* \gamma_2)$ when $B(x) = \gamma_1$.

3.2 Stack Computation Tree Predicate Logic

We are now ready to define our new logic SCTPL. Intuitively, a SCTPL formula is a CTL formula where predicates and regular variable expressions are used as atomic propositions. Using regular variable expressions allows to express predicates on the stack content of the PDS. Moreover, since predicates and regular variable expressions contain variables, we allow quantifiers over variables. For technical reasons, we suppose w.l.o.g. that formulas are given in positive normal form, i.e., negations are applied only to atomic propositions. Indeed, each CTL formula can be written in positive normal form by pushing the negations inside. Moreover, we use the operator \tilde{U} as a dual of the until operator for which the stop condition is not required to occur. Then, standard CTL operators can be expressed as follows: $EF\psi = E[\text{true}U\psi]$, $AF\psi = A[\text{true}U\psi]$, $EG\psi = E[\text{false}\tilde{U}\psi]$ and $AG\psi = A[\text{false}\tilde{U}\psi]$.

More precisely, the set of SCTPL formulas is given by (where $x \in \mathcal{X}$, $a(x_1, \dots, x_n) \in AP_{\mathcal{X}}$ and $e \in \mathcal{R}$):

$$\begin{aligned} \varphi ::= & a(x_1, \dots, x_n) \mid \neg a(x_1, \dots, x_n) \mid e \mid \neg e \mid \varphi \wedge \varphi \mid \varphi \vee \varphi \mid \forall x \varphi \\ & \mid \exists x \varphi \mid AX\varphi \mid EX\varphi \mid A[\varphi U\varphi] \mid E[\varphi U\varphi] \mid A[\varphi \tilde{U}\varphi] \mid E[\varphi \tilde{U}\varphi] \end{aligned}$$

Let φ be a SCTPL formula. The closure $cl(\varphi)$ denotes the set of all the subformulas of φ including φ . The size $|\varphi|$ of φ is the number of elements of $cl(\varphi)$. Let $AP^+(\varphi) = \{a(x_1, \dots, x_n) \in AP_{\mathcal{X}} \mid a(x_1, \dots, x_n) \in cl(\varphi)\}$, $AP^-(\varphi) = \{a(x_1, \dots, x_n) \in AP_{\mathcal{X}} \mid \neg a(x_1, \dots, x_n) \in cl(\varphi)\}$, $Reg^+(\varphi) = \{e \in \mathcal{R} \mid e \in cl(\varphi)\}$, $Reg^-(\varphi) = \{e \in \mathcal{R} \mid \neg e \in cl(\varphi)\}$, and $cl_{\tilde{U}}(\varphi)$ be the set of formulas of $cl(\varphi)$ in the form of $E[\varphi_1 \tilde{U}\varphi_2]$ or $A[\varphi_1 \tilde{U}\varphi_2]$.

Given a PDS $\mathcal{P} = (P, \Gamma, \Delta, \#)$ s.t. $\Gamma \subseteq \mathcal{D}$, let $\lambda : AP_{\mathcal{D}} \rightarrow 2^P$ be a labeling function that assigns a set of control locations to a predicate. Let $c = \langle p, w \rangle$ be a configuration of \mathcal{P} . \mathcal{P} satisfies a SCTPL formula ψ in c , denoted by $c \models_{\lambda} \psi$, iff there exists an environment $B \in \mathcal{B}$ s.t. $c \models_{\lambda}^B \psi$, where $c \models_{\lambda}^B \psi$ is defined by induction as follows:

- $c \models_{\lambda}^B a(x_1, \dots, x_n)$ iff $p \in \lambda(a(B(x_1), \dots, B(x_n)))$.
- $c \models_{\lambda}^B \neg a(x_1, \dots, x_n)$ iff $p \notin \lambda(a(B(x_1), \dots, B(x_n)))$.
- $c \models_{\lambda}^B e$ iff $(c, B) \in L(e)$.
- $c \models_{\lambda}^B \neg e$ iff $(c, B) \notin L(e)$.
- $c \models_{\lambda}^B \psi_1 \wedge \psi_2$ iff $c \models_{\lambda}^B \psi_1$ and $c \models_{\lambda}^B \psi_2$.
- $c \models_{\lambda}^B \psi_1 \vee \psi_2$ iff $c \models_{\lambda}^B \psi_1$ or $c \models_{\lambda}^B \psi_2$.
- $c \models_{\lambda}^B \forall x \psi$ iff $\forall v \in \mathcal{D}$, $c \models_{\lambda}^{B[x \leftarrow v]} \psi$.
- $c \models_{\lambda}^B \exists x \psi$ iff $\exists v \in \mathcal{D}$ s.t. $c \models_{\lambda}^{B[x \leftarrow v]} \psi$.
- $c \models_{\lambda}^B AX \psi$ iff $c' \models_{\lambda}^B \psi$ for every successor c' of c .
- $c \models_{\lambda}^B EX \psi$ iff there exists a successor c' of c s.t. $c' \models_{\lambda}^B \psi$.
- $c \models_{\lambda}^B A[\psi_1 U\psi_2]$ iff for every path $\pi = c_0, c_1, \dots$, of \mathcal{P} with $c_0 = c$, $\exists i \geq 0$ s.t. $c_i \models_{\lambda}^B \psi_2$ and $\forall 0 \leq j < i : c_j \models_{\lambda}^B \psi_1$.

- $c \models_{\lambda}^B E[\psi_1 U \psi_2]$ iff there exists a path $\pi = c_0, c_1, \dots$, of \mathcal{P} with $c_0 = c$ s.t. $\exists i \geq 0$, $c_i \models_{\lambda}^B \psi_2$ and $\forall 0 \leq j < i, c_j \models_{\lambda}^B \psi_1$.
- $c \models_{\lambda}^B A[\psi_1 \tilde{U} \psi_2]$ iff for every path $\pi = c_0, c_1, \dots$, of \mathcal{P} with $c_0 = c$, $\forall i \geq 0$ s.t. $c_i \not\models_{\lambda}^B \psi_2$, $\exists 0 \leq j < i$ s.t. $c_j \models_{\lambda}^B \psi_1$.
- $c \models_{\lambda}^B E[\psi_1 \tilde{U} \psi_2]$ iff there exists a path $\pi = c_0, c_1, \dots$, of \mathcal{P} with $c_0 = c$ s.t. $\forall i \geq 0$ s.t. $c_i \not\models_{\lambda}^B \psi_2$, $\exists 0 \leq j < i$ s.t. $c_j \models_{\lambda}^B \psi_1$.

Intuitively, $c \models_{\lambda}^B \psi$ holds iff the configuration c satisfies the formula ψ under the environment B . Note that a path π satisfies $\psi_1 \tilde{U} \psi_2$ iff either ψ_2 holds everywhere in π , or the first occurrence in the path where ψ_2 does not hold must be preceded by a position where ψ_1 holds.

3.3 Modeling malicious behaviors using SCTPL

SCTPL can be used to precisely specify several malicious behaviors. We needed stack predicates to express most of the specifications. Thus, SCTPL is necessary to specify these behaviors, CTPL is not sufficient. We describe here how e.g., email worms can be specified using SCTPL. The typical behavior of an email worm can be summarized as follows: the worm will first call the API *GetModuleFileNameA* in order to get the name of its executable. For this, the worm needs to call this function with 0 and m as parameters (m corresponds to the address of a memory location), i.e., with $0m$ on the top of the stack since parameters to a function in assembly are passed through the stack. *GetModuleFileNameA* will then write the name of the worm executable on the address m . Then, the worm will copy its file (whose name is at the address m) to other locations using the function *CopyFileA*. It needs to call *CopyFileA* with m as parameter, i.e., with m on the top of the stack. Figure 2(a) shows a disassembled fragment of a code corresponding to this typical behavior. This behavior can be expressed by the SCTPL formula of Figure 2(b). In this formula, Line 2 expresses that there exists a register r_0 such that the address of the memory location m is assigned to r_0 , and such that the value of r_0 does not change until it is pushed onto the stack (subformula $\neg \exists v (mov(r_0, v) \vee lea(r_0, v)) \text{ U } push(r_0)$). Line 3 guarantees that r_0 is not pushed nor popped from the stack until *GetModuleFileNameA* is called, and $0r_0$ is on the top of the stack (the predicate $0r_0 I^*$ ensures this). This guarantees that when *GetModuleFileNameA* is called, r_0 still contains the address of m . Thus, the name of the worm file returned by *GetModuleFileNameA* will be put at the address m . Line 4 is similar to Line 2. It expresses that there exists a register r_1 such that the address of the memory location m is assigned to r_1 , and such that the value of r_1 does not change until it is pushed onto the stack. This guarantees that when r_1 is pushed to the stack, it contains the address of m . Line 5 expresses that r_1 is not pushed nor popped from the stack until *CopyFileA* is called, and r_1 is on the top of the stack (the predicate $r_1 I^*$ ensures this). This guarantees that when *CopyFileA* is called, the value of r_1 is still m . Thus, *CopyFileA* will copy the file whose name is at the address m . Note that we need predicates over the stack to express in a precise manner this specification.

4 SCTPL Model-Checking for Pushdown Systems

In this section, we give an efficient SCTPL model checking algorithm for Pushdown systems. Our procedure works as follows: we reduce this model checking problem to the

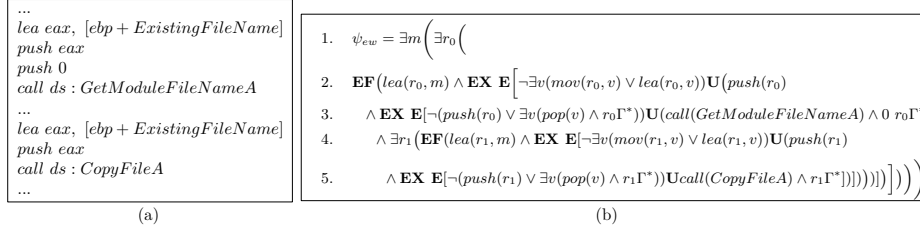


Fig. 2. (a) Email worm (b) Specification of Email worm

emptiness problem in Symbolic Alternating Büchi Pushdown Systems (SABPDS), and we give an algorithm to solve this emptiness problem. To achieve this reduction, we use *variable automata* to represent regular variable expressions. This section is structured as follows. First, we introduce variable automata. Then, we define Symbolic Alternating Büchi Pushdown Systems. Next, we show how SCTPL model checking for PDSs can be reduced to emptiness checking of SABPDSs.

In the remainder of this section, we let \mathcal{X} be a finite set of variables ranging over a finite domain \mathcal{D} , and \mathcal{B} be the set of all the environment functions $B : \mathcal{X} \cup \mathcal{D} \rightarrow \mathcal{D}$.

4.1 Variable Automata

Given a PDS $\mathcal{P} = (P, \Gamma, \Delta, \#)$ s.t. $\Gamma \subseteq \mathcal{D}$, a *Variable Automaton* (VA) is a tuple $M = (Q, \Gamma, \delta, q_0, A)$, where Q is a finite set of states; Γ is the input alphabet; $q_0 \in Q$ is an initial state; $A \subseteq Q$ is a finite set of accepting states; and δ is a finite set of transition rules of the form: $p \xrightarrow{\alpha} \{q_1, \dots, q_n\}$ where α can be x , $\neg x$, or γ , for any $x \in \mathcal{X}$ and $\gamma \in \Gamma$.

Let $B \in \mathcal{B}$. A run of VA on a word $\gamma_1, \dots, \gamma_m$ under B is a tree of height m whose root is labelled by the initial state q_0 , and each node at depth k labelled by a state q has h children labelled by p_1, \dots, p_h , respectively, such that: either $q \xrightarrow{\gamma_k} \{p_1, \dots, p_h\} \in \delta$ and $\gamma_k \in \Gamma$; or $q \xrightarrow{x} \{p_1, \dots, p_h\} \in \delta$, $x \in \mathcal{X}$ and $B(x) = \gamma_k$; or $q \xrightarrow{\neg x} \{p_1, \dots, p_h\} \in \delta$, $x \in \mathcal{X}$ and $B(x) \neq \gamma_k$. A branch of the tree is accepting iff the leaf of the branch is an accepting state. A run is accepting iff all its branches are accepting. A word $\omega \in \Gamma^*$ is accepted by a VA under an environment $B \in \mathcal{B}$ iff the VA has an accepting run on the word ω under the environment B . The language of a VA M , denoted by $L(M)$, is a subset of $(P \times \Gamma^*) \times \mathcal{B}$. $(\langle p, \omega \rangle, B) \in L(M)$ iff M accepts the word ω under the environment B . We can show that:

Theorem 1. *VAs are effectively closed under boolean operations.*

Theorem 2. *For every regular expression $e \in \mathcal{R}$, one can effectively compute in polynomial time a VA M such that $L(M) = L(e)$.*

4.2 Symbolic Alternating Büchi Pushdown Systems

Definition 1. A Symbolic Alternating Büchi Pushdown System (SABPDS) is a tuple $\mathcal{BP} = (P, \Gamma, \Delta, F)$, where P is a finite set of control locations; $\Gamma \subseteq \mathcal{D}$ is the stack

alphabet; $F \subseteq P \times 2^{\mathcal{B}}$ is a set of accepting states; Δ is a finite set of transitions of the form $\langle p, \gamma \rangle \xrightarrow{\mathfrak{X}} [\langle p_1, \omega_1 \rangle, \dots, \langle p_n, \omega_n \rangle]$ where $p \in P$, $\gamma \in \Gamma$, for every i , $1 \leq i \leq n$: $p_i \in P$, $\omega_i \in \Gamma^*$, and $\mathfrak{X} : (\mathcal{B})^n \rightarrow 2^{\mathcal{B}}$ is a function that maps a tuple of environments to a set of environments.

A configuration of a SABPDS is a tuple $\langle [p, B], \omega \rangle$, where $p \in P$ is a control location, $B \in \mathcal{B}$ is an environment and $\omega \in \Gamma^*$ is the stack content. $[p, B] \in P \times \mathcal{B}$ is an accepting state iff $\exists [p, \beta] \in F$ s.t. $B \in \beta$. Let $t = \langle p, \gamma \rangle \xrightarrow{\mathfrak{X}} [\langle p_1, \omega_1 \rangle, \dots, \langle p_n, \omega_n \rangle] \in \Delta$ be a transition, n is the width of the transition t . For every $\omega \in \Gamma^*$, $B, B_1, \dots, B_n \in \mathcal{B}$, if $B \in \mathfrak{X}(B_1, \dots, B_n)$, then the configuration $\langle [p, B], \gamma\omega \rangle$ (resp. $\{\langle [p_1, B_1], \omega_1\omega \rangle, \dots, \langle [p_n, B_n], \omega_n\omega \rangle\}$) is an immediate predecessor (resp. immediate successor) of $\{\langle [p_1, B_1], \omega_1\omega \rangle, \dots, \langle [p_n, B_n], \omega_n\omega \rangle\}$ (resp. $\langle [p, B], \gamma\omega \rangle$). A run ρ of \mathcal{BP} from an initial configuration $\langle [p_0, B_0], \omega_0 \rangle$ is a tree in which the root is labeled by $\langle [p_0, B_0], \omega_0 \rangle$, and the other nodes are labeled by elements of $(P \times \mathcal{B}) \times \Gamma^*$. If a node of ρ labeled by $\langle [p, B], \omega \rangle$ has n children labeled by $\langle [p_1, B_1], \omega_1 \rangle, \dots, \langle [p_n, B_n], \omega_n \rangle$, respectively, then, necessarily, $\langle [p, B], \omega \rangle$ is an immediate predecessor of $\{\langle [p_1, B_1], \omega_1 \rangle, \dots, \langle [p_n, B_n], \omega_n \rangle\}$ in \mathcal{BP} .

A path $c_0c_1\dots$ of a run ρ is an *infinite* sequence of configurations where c_0 is the root of ρ and for every $i \geq 0$, c_{i+1} is one of the children of the node c_i in ρ . The path is accepting iff it visits infinitely often configurations with accepting states. A run ρ is accepting iff all its paths are accepting. Note that an accepting run has only *infinite* paths. A configuration c is accepted (or recognized) by \mathcal{BP} iff \mathcal{BP} has an accepting run starting from c . The language of \mathcal{BP} , denoted by $\mathcal{L}(\mathcal{BP})$, is the set of configurations accepted by \mathcal{BP} .

The predecessor functions $Pre_{\mathcal{BP}}$, $Pre_{\mathcal{BP}}^*$ and $Pre_{\mathcal{BP}}^+$: $2^{(P \times \mathcal{B}) \times \Gamma^*} \rightarrow 2^{(P \times \mathcal{B}) \times \Gamma^*}$ are defined as follows: $Pre_{\mathcal{BP}}(C) = \{c \in (P \times \mathcal{B}) \times \Gamma^* \mid \text{some immediate successor of } c \text{ is a subset of } C\}$, $Pre_{\mathcal{BP}}^*$ is the reflexive and transitive closure of $Pre_{\mathcal{BP}}$, $Pre_{\mathcal{BP}} \circ Pre_{\mathcal{BP}}^*$ is denoted by $Pre_{\mathcal{BP}}^+$.

SABPDS vs. ABPDS. An *Alternating Büchi Pushdown System* (ABPDS for short) [21] can be seen as a SABPDS:

Lemma 1. *Given a SABPDS $\mathcal{BP} = (P, \Gamma, \Delta, F)$, one can compute an equivalent ABPDS \mathcal{BP}' that simulates \mathcal{BP} in $O(|\Delta| \cdot |\mathcal{B}|^{k+1})$ time, where k is the maximum of the widths of the transition rules in Δ and $|\mathcal{B}| = |D|^{|\mathcal{X}|}$.*

Symbolic Alternating Multi-Automata. To finitely represent infinite sets of configurations of SABPDSs, we use Symbolic Alternating Multi-Automata.

Let $\mathcal{BP} = (P, \Gamma, \Delta, F)$ be a SABPDS, a *Symbolic Alternating Multi-Automaton* (SAMA) is a tuple $\mathcal{A} = (Q, \Gamma, \delta, I, Q_f)$, where Q is a finite set of states, Γ is the input alphabet, $\delta \subseteq (Q \times \Gamma) \times 2^Q$ is a finite set of transition rules, $I \subseteq P \times 2^{\mathcal{B}}$ is a finite set of initial states, $Q_f \subseteq Q$ is a finite set of final states. An *Alternating Multi-Automaton* (AMA) is a SAMA such that $I \subseteq P \times \{\emptyset\}$.

We define the reflexive and transitive transition relation $\rightarrow_{\delta} \subseteq (Q \times \Gamma^*) \times 2^Q$ as follows: (1) $q \xrightarrow{\epsilon}_{\delta} \{q\}$ for every $q \in Q$, where ϵ is the empty word, (2) if $q \xrightarrow{\gamma}_{\delta}$

$\{q_1, \dots, q_n\} \in \delta$ and $q_i \xrightarrow{\omega} Q_i$ for every $1 \leq i \leq n$, then $q \xrightarrow{\gamma\omega} \bigcup_{i=1}^n Q_i$. The automaton \mathcal{A} recognizes a configuration $\langle [p, B], \omega \rangle$ iff there exist $Q' \subseteq Q_f$ and $\beta \subseteq \mathcal{B}$ s.t. $B \in \beta$, $[p, \beta] \in I$ and $[p, \beta] \xrightarrow{\omega} Q'$. The language of \mathcal{A} , denoted by $L(\mathcal{A})$, is the set of configurations recognized by \mathcal{A} . A set of configurations is regular if it can be recognized by a SAMA. Similarly, AMAs can also be used to recognize (infinite) regular sets of configurations for ABPDSs.

Proposition 1. *Let $\mathcal{A} = (Q, \Gamma, \delta, I, Q_f)$ be a SAMA. Then, deciding whether a configuration $\langle [p, B], \omega \rangle$ is accepted by \mathcal{A} can be done in $O(|Q| \cdot |\delta| \cdot |\omega| + \tau)$ time, where τ denotes the time used to check whether $B \in \beta$ for some $B \in \mathcal{B}, \beta \subseteq \mathcal{B}$.*

Remark 1. The time τ used to check whether $B \in \beta$ depends on the representation of B and β . In particular, if we use BDDs to represent sets of environment functions, checking whether $B \in \beta$ can be done in $\tau = O(\lceil \log |\mathcal{D}| \rceil \cdot |X|)$ [8].

Computing the language of an SABPDS. We can extend the algorithm of [21] that computes an AMA that recognizes the language of an ABPDS to obtain an algorithm that computes the language of an SABPDS. More precisely:

Theorem 3. *Let $\mathcal{BP} = (P, \Gamma, \Delta, F)$ be a SABPDS, then we can compute a SAMA \mathcal{A} that recognizes $\mathcal{L}(\mathcal{BP})$ in $O(|P|^2 \cdot 2^{2|B|} \cdot |\Gamma| \cdot |\Delta| \cdot 2^{5|P| \cdot 2^{2|B|}})$ time.*

Remark 2. Note that another way to compute $\mathcal{L}(\mathcal{BP})$ is to apply Lemma 1 and produce an equivalent ABPDS \mathcal{BP}' that simulates \mathcal{BP} , and then apply the algorithm of [21] to compute an AMA that recognizes $\mathcal{L}(\mathcal{BP}')$. In practice, in the symbolic case (for SABPDS), the sets of environments β 's can be compactly represented using BDDs for example, whereas in the explicit case (for ABPDS), all the environments B 's have to be considered. Thus, the algorithm behind Theorem 3 will behave better in practice. This is confirmed by the experiments we run where, in the majority of cases, this algorithm terminates in few seconds, whereas if we compute an equivalent ABPDS and apply the algorithm of [21], we run out of memory.

Examples of functions \mathfrak{K} . We give some examples of functions \mathfrak{K} that will be used.

$$- \text{equal}(B_1, \dots, B_n) = \begin{cases} \{B_1\} & \text{if } B_i = B_j \text{ for every } 1 \leq i, j \leq n, \text{ or } n = 1 \\ \emptyset & \text{otherwise.} \end{cases}$$

This function checks that all the B_i 's are equal and returns $\{B_1\}$ (which is equal to $\{B_i\}$ for any i) if this is the case and the emptyset otherwise.

$$- \text{meet}_{\{c_1, \dots, c_n\}}^x(B_1, \dots, B_n) = \begin{cases} \text{Abs}_x(B_1) & \text{if } B_i(x) = c_i \text{ and } B_i(y) = B_j(y) \text{ for } y \neq x, \\ & \text{for every } 1 \leq i, j \leq n, \\ \emptyset & \text{otherwise.} \end{cases}$$

This function checks whether $B_i(x) = c_i$ for every $i, 1 \leq i \leq n$, and for every $y \neq x$ and every $i, j, 1 \leq i, j \leq n$ $B_i(y) = B_j(y)$. It returns $\text{Abs}_x(B_1)$ (which is equal to $\text{Abs}_x(B_i)$ for any i) if this is the case and the emptyset otherwise.

$$- \text{join}_c^x(B_1, \dots, B_n) = \begin{cases} \{B_1\} & \text{if } B_i = B_j \text{ and } B_i(x) = c, \text{ for every } 1 \leq i, j \leq n, \\ \emptyset & \text{otherwise.} \end{cases}$$

This function checks whether $B_i(x) = c$ for every i . If this is the case, it returns $\text{equal}(B_1, \dots, B_n)$, otherwise, it returns the emptyset.

$$- \text{join}_c^{-x}(B_1, \dots, B_n) = \begin{cases} \{B_1\} & \text{if } B_i = B_j \text{ and } B_i(x) \neq c, \text{ for every } 1 \leq i, j \leq n, \\ \emptyset & \text{otherwise.} \end{cases}$$

This function checks whether $B_i(x) \neq c$ for every i . If this is the case, it returns $\text{equal}(B_1, \dots, B_n)$, otherwise, it returns the emptyset.

4.3 From SCTPL Model Checking for PDSs to Emptiness of SABPDS

Let $\mathcal{P} = (P, \Gamma, \Delta, \sharp)$, $\lambda : AP_{\mathcal{D}} \rightarrow 2^P$ be a labeling function, and φ be a SCTPL formula. For every configuration $\langle p, \omega \rangle$, our goal is to determine whether $\langle p, \omega \rangle \models_{\lambda} \varphi$, i.e., whether there exists an environment $B \in \mathcal{B}$ s.t. $\langle p, \omega \rangle \models_{\lambda}^B \varphi$. We proceed as follows: we compute a symbolic alternating Büchi pushdown system \mathcal{BP} s.t. $\langle p, \omega \rangle \models_{\lambda}^B \varphi$ iff $\langle \langle p, \varphi \rangle, B \rangle \in \mathcal{L}(\mathcal{BP})$. Then, $\langle p, \omega \rangle \models_{\lambda} \varphi$ iff there exists $B \in \mathcal{B}$ such that $\langle p, \omega \rangle \models_{\lambda}^B \varphi$.

Let $\text{Reg}^+(\varphi) = \{e_1, \dots, e_k\}$ and $\text{Reg}^-(\varphi) = \{e_{k+1}, \dots, e_m\}$ be the two sets of regular variable expressions¹ that occur in φ . As shown in Theorems 2 and 1, for every i , $1 \leq i \leq k$ we can construct VAs $M_{e_i} = (Q_{e_i}, \Gamma, \delta_{e_i}, s_{e_i}, A_{e_i})$ such that $L(M_{e_i}) = L(e_i)$; and for every j , $k < j \leq m$ we can construct VAs $M_{-e_j} = (Q_{-e_j}, \Gamma, \delta_{-e_j}, s_{-e_j}, A_{-e_j})$ such that $L(M_{-e_j}) = (P \times \Gamma^*) \times \mathcal{B} \setminus L(e_j)$. We suppose w.l.o.g. that the states of these automata are distinct. Let \mathcal{M} be the union of all these automata, \mathcal{F} be the union of all the final states of these automata A_{e_i} 's and A_{-e_j} 's and \mathcal{S} be the union of all the states of these automata Q_{e_i} 's and Q_{-e_j} 's.

Let $\mathcal{BP}_{\varphi} = (P', \Gamma, \Delta', F)$ be the SABPDS defined as follows: $P' = P \times \text{cl}(\varphi) \cup \mathcal{S}$; $F = F_1 \cup F_2 \cup F_3 \cup F_4$, where $F_1 = \{[\langle p, a(x_1, \dots, x_n) \rangle, \beta] \mid a(x_1, \dots, x_n) \in AP^+(\varphi) \text{ and } \beta = \{B \in \mathcal{B} \mid p \in \lambda(a(B(x_1), \dots, B(x_n)))\}\}\}$; $F_2 = \{[\langle p, \neg a(x_1, \dots, x_n) \rangle, \beta] \mid \neg a(x_1, \dots, x_n) \in AP^-(\varphi) \text{ and } \beta = \{B \in \mathcal{B} \mid p \notin \lambda(a(B(x_1), \dots, B(x_n)))\}\}\}$; $F_3 = P \times \text{cl}_{\mathcal{U}}(\varphi) \times \{\mathcal{B}\}$; and $F_4 = \mathcal{F} \times \{\mathcal{B}\}$.

Δ' is the smallest set of transition rules that satisfy the following. For every control location $p \in P$, every subformula $\psi \in \text{cl}(\varphi)$, and every $\gamma \in \Gamma$:

1. if $\psi = a(x_1, \dots, x_n)$ or $\psi = \neg a(x_1, \dots, x_n)$; $\langle \langle p, \psi \rangle, \gamma \rangle \xrightarrow{\text{equal}} \langle \langle p, \psi \rangle, \gamma \rangle \in \Delta'$;
2. if $\psi = \psi_1 \wedge \psi_2$; $\langle \langle p, \psi \rangle, \gamma \rangle \xrightarrow{\text{equal}} [\langle \langle p, \psi_1 \rangle, \gamma \rangle, \langle \langle p, \psi_2 \rangle, \gamma \rangle] \in \Delta'$;
3. if $\psi = \psi_1 \vee \psi_2$; $\langle \langle p, \psi \rangle, \gamma \rangle \xrightarrow{\text{equal}} \langle \langle p, \psi_1 \rangle, \gamma \rangle \in \Delta'$ and $\langle \langle p, \psi \rangle, \gamma \rangle \xrightarrow{\text{equal}} \langle \langle p, \psi_2 \rangle, \gamma \rangle \in \Delta'$;
4. if $\psi = \exists x \psi_1$; $\langle \langle p, \psi \rangle, \gamma \rangle \xrightarrow{\text{meet}_{|c|}^x} \langle \langle p, \psi_1 \rangle, \gamma \rangle \in \Delta'$, for every $c \in \mathcal{D}$;
5. if $\psi = \forall x \psi_1$; $\langle \langle p, \psi \rangle, \gamma \rangle \xrightarrow{\text{meet}_{\mathcal{D}}^x} [\langle \langle p, \psi_1 \rangle, \gamma \rangle, \dots, \langle \langle p, \psi_1 \rangle, \gamma \rangle] \in \Delta'$, where $\langle \langle p, \psi_1 \rangle, \gamma \rangle$ is repeated m times in $[\langle \langle p, \psi_1 \rangle, \gamma \rangle, \dots, \langle \langle p, \psi_1 \rangle, \gamma \rangle]$, where m is the number of elements in \mathcal{D} ;
6. if $\psi = EX\psi_1$; $\langle \langle p, \psi \rangle, \gamma \rangle \xrightarrow{\text{equal}} \langle \langle p', \psi_1 \rangle, \omega \rangle \in \Delta'$ for every $\langle p, \gamma \rangle \hookrightarrow \langle p', \omega \rangle \in \Delta$;
7. if $\psi = AX\psi_1$; $\langle \langle p, \psi \rangle, \gamma \rangle \xrightarrow{\text{equal}} [\langle \langle p_1, \psi_1 \rangle, \omega_1 \rangle, \dots, \langle \langle p_l, \psi_1 \rangle, \omega_l \rangle] \in \Delta'$ such that for every i , $1 \leq i \leq l$, $\langle p, \gamma \rangle \hookrightarrow \langle p_i, \omega_i \rangle \in \Delta$ and these transitions are all the transitions of Δ that have $\langle p, \gamma \rangle$ as left hand side;

¹ $AP^+(\varphi)$, $AP^-(\varphi)$, $\text{Reg}^+(\varphi)$ and $\text{Reg}^-(\varphi)$ are as defined in Section 3.2.

8. if $\psi = E[\psi_1 U \psi_2]$; $\langle (p, \psi), \gamma \rangle \xrightarrow{equal} [\langle (p, \psi_1), \gamma \rangle, \langle (p', \psi), \omega \rangle] \in \mathcal{A}'$ for every rule $\langle p, \gamma \rangle \hookrightarrow \langle p', \omega \rangle \in \mathcal{A}$, and $\langle (p, \psi), \gamma \rangle \xrightarrow{equal} \langle (p, \psi_2), \gamma \rangle \in \mathcal{A}'$;
9. if $\psi = A[\psi_1 U \psi_2]$; $\langle (p, \psi), \gamma \rangle \xrightarrow{equal} [\langle (p, \psi_1), \gamma \rangle, \langle (p_1, \psi), \omega_1 \rangle, \dots, \langle (p_l, \psi), \omega_l \rangle] \in \mathcal{A}'$ such that for every i , $1 \leq i \leq l$, $\langle p, \gamma \rangle \hookrightarrow \langle p_i, \omega_i \rangle \in \mathcal{A}$ and these transitions are all the transitions of \mathcal{A} that have $\langle p, \gamma \rangle$ as left hand side, and $\langle (p, \psi), \gamma \rangle \xrightarrow{equal} \langle (p, \psi_2), \gamma \rangle \in \mathcal{A}'$;
10. if $\psi = E[\psi_1 \tilde{U} \psi_2]$; $\langle (p, \psi), \gamma \rangle \xrightarrow{equal} [\langle (p, \psi_2), \gamma \rangle, \langle (p', \psi), \omega \rangle] \in \mathcal{A}'$ for every $\langle p, \gamma \rangle \hookrightarrow \langle p', \omega \rangle \in \mathcal{A}$, and $\langle (p, \psi), \gamma \rangle \xrightarrow{equal} [\langle (p, \psi_2), \gamma \rangle, \langle (p, \psi_1), \gamma \rangle] \in \mathcal{A}'$;
11. if $\psi = A[\psi_1 \tilde{U} \psi_2]$; $\langle (p, \psi), \gamma \rangle \xrightarrow{equal} [\langle (p, \psi_2), \gamma \rangle, \langle (p_1, \psi), \omega_1 \rangle, \dots, \langle (p_l, \psi), \omega_l \rangle] \in \mathcal{A}'$ such that for every i , $1 \leq i \leq l$, $\langle p, \gamma \rangle \hookrightarrow \langle p_i, \omega_i \rangle \in \mathcal{A}$ and these transitions are all the transitions of \mathcal{A} that have $\langle p, \gamma \rangle$ as left hand side, and $\langle (p, \psi), \gamma \rangle \xrightarrow{equal} [\langle (p, \psi_1), \gamma \rangle, \langle (p, \psi_2), \gamma \rangle] \in \mathcal{A}'$;
12. if $\psi = e$: $\langle (p, \psi), \gamma \rangle \xrightarrow{equal} \langle s_e, \gamma \rangle \in \mathcal{A}'$, where s_e is the initial state of M_e ,
13. if $\psi = \neg e$: $\langle (p, \psi), \gamma \rangle \xrightarrow{equal} \langle s_{\neg e}, \gamma \rangle \in \mathcal{A}'$, where $s_{\neg e}$ is the initial state of $M_{\neg e}$,
14. for every transition $q \xrightarrow{\alpha} \{q_1, \dots, q_n\}$ in \mathcal{M} ; $\langle q, \gamma \rangle \xrightarrow{\mathfrak{R}} \{\langle q_1, \epsilon \rangle, \dots, \langle q_n, \epsilon \rangle\} \in \mathcal{A}'$, where
 - (a) $\mathfrak{R} = equal$ if $\alpha = \gamma$,
 - (b) $\mathfrak{R} = join_\gamma^x$ if $\alpha = x \in \mathcal{X}$,
 - (c) $\mathfrak{R} = join_\gamma^{\neg x}$ if $\alpha = \neg x$ and $x \in \mathcal{X}$,
15. for every $q \in \mathcal{F}$; $\langle q, \# \rangle \xrightarrow{equal} \langle q, \# \rangle \in \mathcal{A}'$.

Roughly speaking, \mathcal{BP}_φ could be seen as the product of \mathcal{P} and φ . \mathcal{BP}_φ recognizes all the configurations $\langle [(p, \psi), B], \omega \rangle$ s.t. $\langle p, \omega \rangle$ satisfies ψ under B . Thus \mathcal{BP}_φ has an accepting run from $\langle [(p, \psi), B], \omega \rangle$ if and only if the configuration $\langle p, \omega \rangle$ satisfies ψ under B . Due to lack of space, we only explain the case $\psi = e$. In this case, the SABPDS \mathcal{BP}_φ accepts $\langle [(p, \psi), B], \omega \rangle$ iff $(\langle p, \omega \rangle, B) \in L(M_e)$. To check this, \mathcal{BP}_φ first goes to state $[s_e, B]$ by Item 12, where s_e is the initial state of M_e , then it continues to check whether ω is accepted by M_e under the environment B . This is ensured by Items 14. Item 14 allows \mathcal{BP}_φ to mimic a run of M_e on ω under the environment B : if \mathcal{BP}_φ is in state $[q, B]$ and the topmost symbol of its stack is γ , then:

- Item 14(a) deals with the case where $q \xrightarrow{\gamma} \{q_1, \dots, q_2\}$ is a transition in δ_e . In this case, \mathcal{BP}_φ moves to the next states $[q_1, B], \dots, [q_n, B]$ while popping γ from the stack. Popping γ allows \mathcal{BP}_φ to check the rest of the word. The function *equal* guarantees that all the environments are the same.
- Item 14(b) deals with the case where $q \xrightarrow{x} \{q_1, \dots, q_2\}$, $x \in \mathcal{X}$ is a transition in δ_e . In this case, \mathcal{BP}_φ can continue to mimic a run of M_e under the environment B only if $B(x) = \gamma$. If this holds, \mathcal{BP}_φ moves to the next states $[q_1, B], \dots, [q_n, B]$ and pops γ from the stack, which allows \mathcal{BP}_φ to check the rest content of the stack. The function *join $_\gamma^x$* ensures that all the environments are the same and the value of $B(x)$ is γ .
- Similarly, Item 14(c) deals with the case where $q \xrightarrow{\neg x} \{q_1, \dots, q_2\}$ is in δ_e .

Thus, $(\langle p, \omega \rangle, B) \in L(M_e)$ iff M_e reaches final states f_1, \dots, f_n of M_e after reading the word ω , i.e., iff \mathcal{BP}_φ reaches a set of states $[f_1, B], \dots, [f_n, B]$ with an empty stack (a stack containing only the bottom stack symbol $\#$). This is why F_4 is a set of accepting states. Moreover, since all the accepting paths are infinite, Item 15 adds a loop on every

configuration $\langle [f, B], \sharp \rangle$ where f is a final state of M and \sharp is the stack symbol (this makes the paths of \mathcal{BP}_φ that reach a state $\langle [f, B], \sharp \rangle$ accepting). Formally, we can show:

Theorem 4. *Given a PDS $\mathcal{P} = (P, \Gamma, \Delta, \sharp)$, a function $\lambda : \text{AP}_{\mathcal{D}} \rightarrow 2^P$, a SCTPL formula φ , and a configuration $\langle p, \omega \rangle$ of \mathcal{P} , we have: for every $B \in \mathcal{B}$, $\langle p, \omega \rangle \models_\lambda^B \varphi$ iff \mathcal{BP}_φ has an accepting run from the configuration $\langle [p, \varphi], B, \omega \rangle$.*

4.4 SCTPL model-checking for PDSs

Given a PDS $\mathcal{P} = (P, \Gamma, \Delta, \sharp)$, a labeling function λ , and a SCTPL formula φ , thanks to Theorems 4 and 3, and due to the fact that \mathcal{BP}_φ has $O(|P| \cdot |\varphi| + k)$ states and $O((|P| \cdot |\Gamma| + |\Delta|) \cdot |\varphi| + d)$ transitions, where k and d are the number of states and the number of transitions of the union \mathcal{M} of the Variable Automata involved in φ ; we get the following:

Corollary 1. *Given a PDS $\mathcal{P} = (P, \Gamma, \Delta, \sharp)$, a SCTPL formula φ and a labeling function λ , we can effectively compute a SAMA \mathcal{A} in time $O((|P||\varphi| + k)^2 \cdot 2^{2^{|\mathcal{B}|}} \cdot |\Gamma| \cdot (|P||\Gamma| + |\Delta|)|\varphi| + d) \cdot 2^{5(|P||\varphi| + k) \cdot 2^{2^{|\mathcal{B}|}}}$, where k is the number of states of \mathcal{M} and d is the number of transition rules of \mathcal{M} such that for every configuration $\langle p, \omega \rangle$ of \mathcal{P} :*

1. $\langle p, \omega \rangle \models_\lambda \varphi$ iff there exists a $B \in \mathcal{B}$ s.t. \mathcal{A} recognizes $\langle [p, \varphi], B, \omega \rangle$.
2. for every $B \in \mathcal{B}$: $\langle p, \omega \rangle \models_\lambda^B \varphi$ iff \mathcal{A} recognizes $\langle [p, \varphi], B, \omega \rangle$.

Thus, thanks to this corollary and to Proposition 1, it follows that it is possible to determine whether a PDS configuration satisfies a SCTPL formula:

Corollary 2. *It is possible to decide whether a PDS configuration satisfies a SCTPL formula.*

Remark 3. We can transform every SCTPL formula ψ to an equivalent CTL with regular valuations formula ψ' . Then, applying [21], we can construct an AMA recognizing all the configurations which satisfy ψ' . However, in practice, thanks to the compact representation of the sets of environments β 's using BDDs, model-checking SCTPL using our symbolic techniques behaves much better than reducing SCTPL to CTL with regular valuations and then applying [21]. Indeed, the experiments we run show that in most of the cases, our symbolic algorithm for SCTPL model-checking terminates in few seconds, whereas translating the SCTPL formula to CTL with regular valuations and then applying [21] would run out of memory.

5 Experiments

We implemented our techniques in a tool for malware detection. We use IDAPro [3] as disassembler. We use BDDs to represent sets of environments. We carried out different experiments. We obtained interesting results. In particular, our tool was able to detect several viruses taken from [14]. Our results are reported in Table 1. **Column** $|P|$ gives the number of control locations of the PDS model. Every program is checked against several malicious behaviors. A program is declared as a potential virus if it satisfies one

Examples	P	Our techniques		SABPDS→ABPDS		SCTPL→CTLr		Result
		Time(s)	Mem(Mb)	Time(s)	Mem(Mb)	Time(s)	Mem(Mb)	
Klez.a	42	1.62	10.8	-	MemOut	-	MemOut	Y
Klez.b	45	1.55	10.8	-	MemOut	-	MemOut	Y
Klez.c	41	1.27	8.9	-	MemOut	-	MemOut	Y
Klez.d	51	1.47	10.3	-	MemOut	-	MemOut	Y
Klez.e	52	0.77	7.0	-	MemOut	-	MemOut	Y
Klez.f	50	0.76	7.0	-	MemOut	-	MemOut	Y
Klez.g	47	0.75	7.0	-	MemOut	-	MemOut	Y
Klez.i	49	0.74	7.0	-	MemOut	-	MemOut	Y
Klez.j	55	0.74	7.0	-	MemOut	-	MemOut	Y
Mydoom.c	210	145.20	322.8	-	MemOut	-	MemOut	Y
Mydoom.e	288	123.22	267.5	-	MemOut	-	MemOut	Y
Mydoom.g	256	117.50	256.7	-	MemOut	-	MemOut	Y
Predec.j	25	0.23	0.81	-	MemOut	56.14	36.16	Y
Netsky.a	69	2.73	14.5	-	MemOut	-	MemOut	Y
Akez	42	0.22	0.3	-	MemOut	0.44	2.49	Y
Netsky.b	80	2.73	14.5	-	MemOut	-	MemOut	Y
Netsky.c	78	2.73	14.5	-	MemOut	-	MemOut	Y
Netsky.d	72	2.73	14.5	-	MemOut	-	MemOut	Y
Alcau.h	48	0.83	0.9	-	MemOut	1.14	6.88	Y
Uedit32	180	92.58	100.94	-	MemOut	-	MemOut	N
Alcau.l	2	0.30	0.7	-	MemOut	0.86	3.96	Y
Cygwin32	212	23.72	123.31	-	MemOut	-	MemOut	N
cmd.exe	202	1.44	25.52	-	MemOut	-	MemOut	N
Alcau.o	68	0.20	0.6	-	MemOut	0.83	3.37	Y
Mydoor.a	256	113.2	227.4	-	MemOut	-	MemOut	Y

Examples	P	Our techniques		SABPDS→ABPDS		SCTPL→CTLr		Result
		Time(s)	Mem(Mb)	Time(s)	Mem(Mb)	Time(s)	Mem(Mb)	
Adson.1559	52	0.22	2.1	-	MemOut	-	MemOut	Y
Adson.1651	54	0.23	2.1	-	MemOut	-	MemOut	Y
Adson.1703	55	0.25	2.1	-	MemOut	-	MemOut	Y
Adson.1734	54	0.31	2.6	-	MemOut	-	MemOut	Y
Alcaul.d	62	0.20	0.8	-	MemOut	47.70	51	Y
Alcaul.i	88	4.38	0.28	-	MemOut	159.88	169.64	Y
Alcaul.j	79	0.30	2.1	-	MemOut	218.25	198.71	Y
Oroch.3982	89	3.70	7.72	-	MemOut	-	MemOut	Y
KME	145	999.31	20.04	-	MemOut	-	MemOut	Y
Anar.a	41	1.16	1.60	885.33	343.24	54.92	34.12	Y
Anar.b	47	1.49	1.60	891.42	348.54	56.14	36.16	Y
Atak.b	126	762.34	18.15	-	MemOut	-	MemOut	Y
Alcaul.c	33	0.12	0.3	-	MemOut	0.41	2.19	Y
Bagle.d	88	652.23	16.96	-	MemOut	-	MemOut	Y
Alcaul.f	52	0.09	0.3	-	MemOut	0.53	2.23	Y
Alcaul.b	50	0.06	0.2	-	MemOut	0.28	1.18	Y
Alcaul.e	49	0.49	0.9	-	MemOut	1.03	5.28	Y
Alcaul.g	53	0.31	0.7	-	MemOut	0.97	4.45	Y
Evo.l.a	102	9.58	3.22	-	MemOut	-	MemOut	Y
Alcaul.k	52	0.26	0.6	-	MemOut	0.76	3.65	Y
Alcaul.m	53	0.20	0.6	-	MemOut	0.88	3.37	Y
Alcaul.n	34	0.12	0.3	-	MemOut	0.44	2.28	Y
Klinge	78	237.50	4.49	-	MemOut	0.83	3.37	Y
Atak.f	220	23.4	139.1	-	MemOut	-	MemOut	Y
Mydoor.a.y	328	124.2	232.5	-	MemOut	-	MemOut	Y

Table 1. Detection of real malwares

of the specifications. **Column** *time(s)* and *mem(Mb)* give the time (in seconds) and the memory (in Mb). The last **Column** *result* is *Y* if the program contains the malicious behaviors given in **Column** *Formula*, and *N* if not. We also compared our techniques against translating SABPDS to ABPDS (**Columns** “SABPDS→ABPDS”), or translating SCTPL to CTL with regular valuations (**Columns** “SCTPL→CTLr”). We were able to detect all the viruses that we considered, whereas applying the translation from S-ABPDS to ABPDS or from SCTPL to CTL with regular valuations would run out of memory in most of the cases, and thus cannot detect the viruses. Our tool was also able to deduce that some benign programs are not viruses. E.g. we tried the following benign programs: *Uedit32*, a fragment of Ultra Edit Text Editor software by IDM Computer Solutions; *Cygwin32* a fragment of the Setup software of Cygwin, a Linux-like environment for Windows. *cmd.exe* is the Microsoft-supplied command-line interpreter.

Moreover, we run several experiments to check how robust are our techniques in virus detection in case the virus writers use obfuscation techniques. To this aim, we considered some of the viruses of Table 1, and we added several obfuscations manually such as: instruction reordering (reordering the instructions inside the code and using jump instructions so that the control flow is not changed), dead code insertion, register renaming, splitting the code into several procedures, adding useless stack operations, etc. We tested 5 variants for each type of obfuscation of the viruses Mydoom.g, Netsky.a, Bagle.d, Adson.1734 and Akez. The results are reported in Table 2. Our techniques were able to detect all these variations, whereas the three well known and widely used free antiviruses *Avira* [2], *Qihoo 360* [4] and *Avast* [1] were not able to detect several of these virus variations.

References

1. Avast antivirus, free version. <http://www.avast.com>.
2. Avira antivirus, free version. <http://www.avira.com>.
3. IDA Pro. <http://www.hex-rays.com/idapro/>.

Obfuscation	Our techniques detection rate	Avira antivirus detection rate	Qihoo 360 antivirus detection rate	Avast antivirus detection rate
nop-insertion	100%	65%	55%	60%
code-reordering	100%	40%	35%	45%
register-renaming	100%	25%	25%	30%
stack-operation	100%	20%	25%	20%
procedure-split	100%	5%	5%	5%

Table 2. Detection of obfuscated Viruses

4. Qihoo 360 antivirus. <http://www.360.cn>.
5. G. Balakrishnan, T. W. Reps, N. Kidd, A. Lal, J. Lim, D. Melski, R. Gruian, S. H. Yong, C.-H. Chen, and T. Teitelbaum. Model checking x86 executables with codesurfer/x86 and wpds++. In *CAV*, 2005.
6. J. Bergeron, M. Debbabi, J. Desharnais, M. Erhioui, Y. Lavoie, and N. Tawbi. Static detection of malicious code in executable programs. In *SREIS*, 2001.
7. G. Bonfante, M. Kaczmarek, and J.-Y. Marion. Architecture of a Morphological Malware Detector. *Journal in Computer Virology*, 5:263–270, 2009.
8. R. E. Bryant. Symbolic boolean manipulation with ordered binary-decision diagrams. *ACM Comput. Surv.*, 24(3), 1992.
9. M. Christodorescu and S. Jha. Static analysis of executables to detect malicious patterns. In *12th USENIX Security Symposium*, 2003.
10. M. Christodorescu, S. Jha, and C. Kruegel. Mining specifications of malicious behavior. In *ISEC*, 2008.
11. M. Christodorescu, S. Jha, S. A. Seshia, D. X. Song, and R. E. Bryant. Semantics-aware malware detection. In *IEEE Symposium on Security and Privacy*, 2005.
12. J. Esparza, A. Kucera, and S. Schwoon. Model checking LTL with regular valuations for pushdown systems. *Inf. Comput.*, 186(2), 2003.
13. J. Esparza and S. Schwoon. A bdd-based model checker for recursive programs. In *CAV*, 2001.
14. V. Heavens. <http://vx.netlux.org>.
15. A. Holzer, J. Kinder, and H. Veith. Using verification technology to specify and detect malware. In *EUROCAST*, 2007.
16. J. Kinder, S. Katzenbeisser, C. Schallhart, and H. Veith. Detecting malicious code by model checking. In *DIMVA*, 2005.
17. J. Kinder, S. Katzenbeisser, C. Schallhart, and H. Veith. Proactive detection of computer worms using model checking. *IEEE Transactions on Dependable and Secure Computing*, 7(4), 2010.
18. A. Lakhota, D. R. Boccoardo, A. Singh, and A. Manacero. Context-sensitive analysis of obfuscated x86 executables. In *PEPM*, 2010.
19. A. Lakhota, E. U. Kumar, and M. Venable. A method for detecting obfuscated calls in malicious binaries. *IEEE Trans. Software Eng.*, 31(11), 2005.
20. P. K. Singh and A. Lakhota. Static verification of worm and virus behavior in binary executables using model checking. In *IAW*, 2003.
21. F. Song and T. Touili. Efficient CTL model-checking for pushdown systems. In *CONCUR*, 2011.