

## Behavior Abstraction in Malware Analysis

Philippe Beaucamps, Isabelle Gnaedig, Jean-Yves Marion

► **To cite this version:**

Philippe Beaucamps, Isabelle Gnaedig, Jean-Yves Marion. Behavior Abstraction in Malware Analysis. 1st International Conference on Runtime Verification, Howard Barringer, Klaus Havelund, Insup Lee, Nov 2010, St. Julians, Malta. pp.168-182, 10.1007/978-3-642-16612-9\_14 . inria-00536500

**HAL Id: inria-00536500**

**<https://hal.inria.fr/inria-00536500>**

Submitted on 10 Oct 2018

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

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

# Behavior Abstraction in Malware Analysis\*

Philippe Beaucamps, Isabelle Gnaedig, Jean-Yves Marion

INPL - INRIA Nancy Grand Est - Nancy-Université - LORIA  
Campus Scientifique - BP 239 F54506 Vandoeuvre-lès-Nancy Cedex, France

**Email:** {Philippe.Beaucamps, Isabelle.Gnaedig, Jean-Yves.Marion}@loria.fr

**Abstract.** We present an approach for proactive malware detection working by abstraction of program behaviors. Our technique consists in abstracting program traces, by rewriting given subtraces into abstract symbols representing their functionality. Traces are captured dynamically by code instrumentation, which allows us to handle packed or self-modifying malware. Suspicious behaviors are detected by comparing trace abstractions to reference malicious behaviors. The expressive power of abstraction allows us to handle general suspicious behaviors rather than specific malware code and then, to detect malware mutations. We present and discuss an implementation validating our approach.

**Keywords:** Malware, behavioral detection, behavior abstraction, trace, string rewriting, finite state automaton, formal language, dynamic binary instrumentation.

## 1 Introduction

Detection techniques of computer malware have traditionally relied on a combination of static and dynamic analysis. A shortcoming of static analysis, however, is the general intractability of knowing in advance the entire program code as it may change dynamically. Packing and obfuscation techniques typically capitalize on this intractability to prevent the reconstruction of the program code. Structural and behavioral techniques, on the other hand, may be used to guard against code protection and code transformation and are consequently more robust and reliable. These techniques rely on the analysis of the program structure or its behavior rather than its binary code. This ensures an independence from basic syntactic modifications or from packing techniques.

Structural techniques analyze the control flow graph representing the program, by assuming that its structure remains untouched. They compare it to control flow graphs of known malware using static analysis [10], sub-graph isomorphism [9], tree automata [7] or similarity measures [12,4]. Unfortunately,

---

\* This work has been partially supported by the High Security Lab of the LORIA in Nancy: <http://lhs.loria.fr>. A full version of this article can be found in [5].—— The original publication is available at [www.springerlink.com](http://www.springerlink.com) —— [https://link.springer.com/chapter/10.1007/978-3-642-16612-9\\_14](https://link.springer.com/chapter/10.1007/978-3-642-16612-9_14)

these techniques are not resilient to functional polymorphism of behaviors, of which malware variants are a form, and which expresses that behaviors can be carried out in different ways, without their functionality being altered. This polymorphism often impacts the structure of the control flow graph.

Conversely, behavioral approaches, first introduced in Cohen’s seminal work [11], monitor system calls and their arguments and have traditionally relied on the use of finite state machines [18,21]. Recent approaches [19] deal with functional polymorphism by preprocessing execution traces and transforming them into a high-level representation which captures their semantic meaning. But as these approaches deal with the execution trace being observed, they analyze a single behavior at a time. Subsequently, [16] proposed to use attribute automata but the cost is an exponential time complexity procedure.

Other behavioral approaches also use model checking techniques to track data [6,17,22]: they define behavioral signatures as temporal logic formulae, defined on a syntactic level. But none of these approaches considers functional polymorphism. Moreover they do not tackle either the problem of constructing a high-level view of a program, which limits their applicability.

Our goal here is to provide a generic framework for malware detection, abstract enough to be independent of the implementation of programs, resilient to variants and relying on general suspicious behaviors rather than on specific malware code. Unlike the approaches cited before, which only consider the detection scenario by working on one trace at a time, we intend to make our formalism more generally applicable to analysis and signature generation of unknown malware by working on a set of traces representing the whole behavior of a program.

For this purpose, we present an approach working on an abstract representation of the behavior of a program. We propose an original strong theoretical setting underpinned by the theory of formal languages and based on string rewriting systems and finite state automata. Abstraction is carried out with respect to behavior patterns defined by string rewriting systems. Behavioral detection is then carried out by automata intersection.

More precisely, execution traces of a program describe the capture of specific data such as program instructions, system calls with their arguments, or file system and network interactions. We represent a set of execution traces by an automaton called trace automaton. An abstraction of this trace automaton is then constructed, with respect to a set of predefined behavior patterns defined as regular languages and describing high-level properties or actions such as a file operation, a hook installation or a data leak. This gives a representation independent of the program implementation. Finally, the abstracted trace automaton is intersected with a malware database, composed of abstract signatures representing known malicious behaviors.

Our technique offers two detection scenarios: identifying a program as suspicious when it performs some malicious action like keylogging, or detecting an action sequence similar to the one of a known malware. The model we use, combining string rewriting systems and finite state automata, allows us to detect very efficiently malicious behaviors with high level descriptions, that is in linear

time. Detection speed can be tuned by setting the right tier of abstraction level. So our behavioral detection model could be used inside a firewall for example.

After presenting the background in Section 2, we define behavior patterns in Section 3. Section 4 presents the abstraction mechanism of trace languages. Section 5 explains how to represent trace languages by trace automata and gives complexity bounds for computing abstractions. Section 6 formalizes the detection problem with its cost. Section 7 presents the implementation of our approach together with experiments. We conclude in Section 8.

## 2 Background

Let  $\Sigma$  be some finite alphabet. We denote by  $\Sigma^*$  the set of finite words on  $\Sigma$ . Subsets of  $\Sigma^*$  are called languages on  $\Sigma$ . The empty word is denoted by  $\epsilon$ .

Let  $\Sigma'$  be some finite alphabet. The projection homomorphism which maps words of  $\Sigma^*$  to words of  $\Sigma'^*$  is denoted by  $u|_{\Sigma'}$ , for any  $u \in \Sigma^*$  and is defined, for  $a \in \Sigma$ , by:  $a|_{\Sigma'} = a$  if  $a \in \Sigma'$  and  $a|_{\Sigma'} = \epsilon$  otherwise. This definition is homeomorphically extended to languages on  $\Sigma$ , and the projection on  $\Sigma'$  of a language  $L$  on  $\Sigma$  is denoted by  $L|_{\Sigma'}$ .

A finite state automaton  $\mathcal{A}$  on an alphabet  $\Sigma$  is a tuple  $(Q, \delta, q_0, F)$  where  $Q$  is a finite set of states,  $\delta : Q \times \Sigma \rightarrow Q$  is the transition relation,  $q_0 \in Q$  is the initial state, and  $F \subseteq Q$  is the set of final states. A run of  $\mathcal{A}$  on a word  $w = a_0a_1 \cdots a_n$  is a sequence of states  $r = q_0q_1 \cdots q_{m \leq n+1}$  such that:  $\forall i < m, q_{i+1} \in \delta(q_i, a_i)$ . The run  $r$  is successful if  $m = n + 1$  and  $q_m \in F$ ; in this case,  $w$  is said to be recognized by  $\mathcal{A}$ . The set of words for which there exists a successful run of  $\mathcal{A}$  is the language recognized by  $\mathcal{A}$  and is denoted by  $\mathcal{L}(\mathcal{A})$ . Languages recognized by some finite state automaton  $\mathcal{A}$  are called regular. The size of a finite state automaton  $A$ , denoted by  $|A|$ , is defined as the number of states of  $A$ .

For a given binary relation  $\rightarrow$ , we denote by  $\rightarrow^*$  its reflexive transitive closure. A string rewriting system (SRS in short) is a triple  $(\Sigma, V, \rightarrow)$ , where  $V$  is a set of variables and  $\rightarrow$  a binary relation on  $(\Sigma \cup V)^*$ . A regular SRS is a 4-tuple  $(\Sigma, V, S, \rightarrow)$ , where  $S \in V$  and the relation  $\rightarrow$  is generated by rules of the form  $a \rightarrow A, aA \rightarrow B$  and  $\epsilon \rightarrow A$ , with  $a \in \Sigma, A, B \in V$ . The language recognized by a regular SRS is the set  $\{u|u \in \Sigma^*, u \rightarrow^* S\}$ . Languages recognized by regular SRS are exactly the regular languages. The size of a regular SRS  $R$ , denoted by  $|R|$ , is defined as the number of variables of  $V$ .

## 3 Behaviors

We now introduce a model of abstract machine from which we define notions of execution trace and behavior. An abstract machine  $\mathcal{M}$  consists of the triple  $(\mu_0, IP_0, \rightarrow)$  where  $(\mu_0, IP_0)$  is an initial configuration of  $\mathcal{M}$  and  $\rightarrow$  is a transition function from *Configurations* to *Configurations*, where *Configurations* denotes the set of configurations of  $\mathcal{M}$ .

A configuration of  $\mathcal{M}$  is a pair  $(\mu, IP)$  where:

- $\mu : \text{Addresses} \rightarrow \text{Data}$  represents the memory of  $\mathcal{M}$ . *Addresses* is the set of addresses of  $\mathcal{M}$  and *Data* is a set of values; both are subsets of  $\mathbb{N}$ ;
- $IP \in \text{Addresses}$  is the instruction pointer.

Thus, we have  $(\mu, IP) \rightarrow (\mu', IP')$  if the machine  $\mathcal{M}$  executes the instruction at address  $IP$  of memory  $\mu$ . The memory  $\mu'$  is the memory obtained after executing this instruction and  $IP'$  is the address of the next instruction to execute. A program is a set of instructions. An *execution* of an abstract machine  $\mathcal{M}$  is a finite sequence:

$$(\mu_0, IP_0) \rightarrow (\mu_1, IP_1) \rightarrow \dots \rightarrow (\mu_n, IP_n).$$

In our scenario, the configuration  $(\mu_0, IP_0)$  is the initial configuration. A program is loaded into  $\mu_0$  at the address pointed by the instruction pointer  $IP_0$ . So, at the beginning of an execution, a program is executed inside an initial environment (also) given by  $\mu_0$ . Then, the program is run. At each step, we see interactions with the "outside" through the memory. Though our model of abstract machine can only represent single-threaded programs, this formalization of communications is enough for our purpose.

The reader will notice that we focus on abstract machines rather than on programming languages. There are several reasons for this. First of all, our model allows us to talk about programming languages at any level of abstraction. Second, in the context of malware, programs are generally self-modifying. Programs are treated as data objects and elements of *Data* are regarded as instructions. A program text is variable, which is not the usual point of view in semantics, see for example Gunter's textbook [15]. Moreover, low level instructions, like in the x86 case, are not of fixed size and a program can modify its code by instruction misalignment. So we think that our model of abstract machine is a right approach to underpin our study on malware behavior detection.

As said before, dynamic analysis and detection rely on capture and analysis of execution data. This data may be the sequence of instructions being executed, the sequence of system calls (i.e. calls to system code), etc. Other approaches may capture yet higher-level actions of a program, for instance network interactions, file system accesses, inter process communications (IPC), register usage statistics or any data that can be used to characterize a behavior. Our framework aims at dealing with any kind of the above data using an alphabet  $\Sigma$ .

We first formalize the capture of some execution data, represented by elements of  $\Sigma$ , in the machine  $\mathcal{M}$ .

**Definition 1 (Capture operator).** *A capture operator with respect to  $\Sigma$  is an operator  $\pi : \text{Configurations} \rightarrow \Sigma \cup \{\epsilon\}$  which associates with some configuration the captured data if any, and  $\epsilon$  otherwise.*

Note that in the general case, captured data of a configuration  $c_n$  may depend on the history of the execution i.e. on configurations  $c_{i_1} \dots c_{i_k}$  for  $i_1 \dots i_k \in [1..n-1]$  at previous times. This is needed for example to compute register usage statistics or to capture non atomic behaviors (for instance "smtp connection" which is the

combination of a network connection and reception of the message “220 .\* SMTP Ready”). For the sake of simplicity, for defining  $\pi$ , we do not consider  $c_{i_1} \dots c_{i_k}$ .

From now on, in our examples, we consider more specifically the operator capturing library calls, including system calls.

**Definition 2 (Execution trace).** *Let  $\mathcal{M}$  be a machine,  $e = c_1 \dots c_n$  an execution of  $\mathcal{M}$  and  $\pi$  a capture operator with respect to  $\Sigma$ . Then  $\pi(c_1) \dots \pi(c_n) \in \Sigma^*$  is the execution trace of the execution  $e$  of  $\mathcal{M}$  with respect to  $\pi$ , denoted  $\pi(e)$ .*

In the following, we will call trace language of a machine  $\mathcal{M}$  with respect to the capture operator  $\pi$ , the set of execution traces of  $\mathcal{M}$  with respect to  $\pi$ . We denote it by  $Traces_\pi(\mathcal{M})$ , or simply  $Traces(\mathcal{M})$  when  $\pi$  is clear from the context. We can now formally define the notion of behavior we want to detect in a program, with respect to some capture operator.

**Definition 3 (Behavior pattern).** *A regular behavior pattern  $\mathcal{B}$  is a regular language on  $\Sigma$ . We call simple behavior pattern an element of  $\mathcal{B}$ .*

A behavior pattern does not necessarily describe a malicious behavior in itself. It can describe an innocuous behavior – e.g. creating a process, opening a file, sending a mail – or a relevant behavior sequence, possibly modulo shuffling of independent actions. It can also represent a more specific behavior which is shared by different malicious codes.

*Example 1.* Throughout this paper, we consider the example of the Allapple worm, a polymorphic network worm. A simplified excerpt of the code of its variant Allapple.A is given in [5]. It contains three behavior patterns: the ping of a remote host, the opening of a Netbios connection and the scanning of local drives. An example of execution trace of this excerpt is the following sequence of library calls:

```
...GetLogicalDriveStrings.GetDriveType.FindFirstFile.FindFirstFile.  
FindNextFile...
```

This trace exhibits the behavior pattern which describes the scanning of local drives: `GetLogicalDriveStrings.GetDriveType.FindFirstFile`.

## 4 Trace Abstraction

Given some machine, recall that our goal is to abstract its execution traces, by replacing concrete behavior patterns by more abstract representations, expressing their functionality, in order to compare them to reference malicious behaviors.

In this section, we formally define the abstraction of a program trace with respect to behavior patterns. We start with a simple behavior pattern, then generalize abstraction to a regular behavior pattern and finally to a set of regular behavior patterns. We show that abstraction can be viewed as a string rewriting process where each behavior pattern is replaced by a new specific symbol.

#### 4.1 Abstracting with respect to a Simple Behavior Pattern

Let  $t \in \mathcal{B}$  be a simple behavior pattern for some  $\mathcal{B}$ . Let  $\mathcal{M}$  be a machine and  $u \in \text{Traces}_\pi(\mathcal{M})$  an execution trace of  $\mathcal{M}$  for some  $\pi$ . Identifying occurrences of the pattern  $t$  in the trace  $u$  amounts to matching  $t$  to some subword of  $u$ .

Let  $\lambda \notin \Sigma$  be a new symbol denoting the abstraction of our pattern  $t$ . We define  $\Sigma' = \Sigma \cup \{\lambda\}$ . Trace abstraction of  $u$  with respect to pattern  $t$  of  $\Sigma^*$  is defined by rewriting  $u$  with the string rewriting system  $R$  on  $\Sigma'^*$ , where  $R$  is composed of a single rewrite rule:  $R = \{t \rightarrow \lambda\}$ .

Let  $\rightarrow_t$  denote the rewriting relation induced by  $R$ , which rewrites substrings:

$$\forall u, v \in \Sigma'^*, u \text{ rewrites into } v, \text{ which is written } u \rightarrow_t v \text{ iff}$$

$$\exists u' \in \Sigma'^*, \exists u'' \in \Sigma'^*, \begin{cases} u = u' \cdot t \cdot u'' \\ v = u' \cdot \lambda \cdot u'' \end{cases}$$

A trace has been abstracted when every occurrence of the pattern  $t$  has been replaced by the abstract symbol  $\lambda$ . Thus abstraction of a trace is nothing but normalization of that trace with respect to the SRS  $R$ . Note that  $R$  is not confluent in general, so a trace may have several normal forms. Abstraction of a trace can be naturally generalized to a trace language.

*Example 2.* Returning to our excerpt of the Allapple worm, suppose we are interested in detecting the previously defined behavior of scanning local drives. Then  $t = \text{GetLogicalDriveStrings.GetDriveType.FindFirstFile}$ , and we define  $\lambda = \text{SCAN\_DRIVES}$ . The execution trace in Example 1 is thus abstracted into:  
`...SCAN_DRIVES.FindFirstFile.FindNextFile...`

Note that we consider normal forms instead of partially reduced ones. In other words, we require that every occurrence of a pattern that can be matched is eventually rewritten. This allows to ensure that the computation of the malicious behavior exhibited by a given malware is maximal. Also, this allows a more precise detection of behaviors of the type ‘‘A followed by B, without C in between’’: partially rewriting a trace  $acb$  into  $AcB$  would then lead to a positive detection, while the normal form  $ACB$  does not match the behavior.

#### 4.2 Abstracting with respect to a Regular Behavior Pattern

A behavior pattern can usually be achieved in different ways and malware writers may capitalize on this to perform functional polymorphism. For instance, creating a file may be processed using different sequences of instructions: detection will only be interested in the functionality expressed by these sequences, not in their implementation details. Thus, we now consider the case of a behavior pattern defined as a language, and more specifically as a regular language, as presented in Definition 3. Indeed, regular patterns allow to account for example for behaviors interleaved with other unrelated or irrelevant actions. For instance, file writes may often be interleaved with other interactions with the system, e.g. string manipulation functions.

Normalization now consists in rewriting a trace language with the rules transforming into  $\lambda$  any simple pattern of a behavior pattern  $\mathcal{B}$ , where  $\mathcal{B}$  is a regular language on  $\Sigma^*$ . Since  $\mathcal{B}$  may be infinite, we consider a regular SRS  $(\Sigma, V, S, \rightarrow)$  recognizing  $\mathcal{B}$  and we define relation  $\rightarrow_{\mathcal{B}}$  on  $(\Sigma' \cup V)^*$  as the rewriting relation induced by  $\rightarrow \cup \{S \rightarrow \lambda\}$ . Using the reflexive transitive closure of  $\rightarrow_{\mathcal{B}}$ , we define trace language abstraction with respect to a regular behavior pattern.

**Definition 4 (Abstract trace language).** *Let  $\mathcal{B}$  be a regular behavior pattern, and  $\rightarrow_{\mathcal{B}}$  its associated rewriting relation. The abstract form of the trace language  $L$  with respect to  $\mathcal{B}$ , denoted by  $L \downarrow_{\mathcal{B}}$ , is defined by:*

$$L \downarrow_{\mathcal{B}} = \{v \in \Sigma'^* \mid \exists u \in L, u \rightarrow_{\mathcal{B}}^* v \text{ and } \nexists w \in \Sigma'^*, v \rightarrow_{\mathcal{B}} w\}.$$

The following theorem expresses that abstraction preserves regularity, which is a fundamental characteristic of our approach.

**Theorem 1.** *Let  $\mathcal{B}$  be a regular behavior pattern and  $L$  a trace language. If  $L$  is regular, then so is  $L \downarrow_{\mathcal{B}}$ .*

Proofs of all theorems of the paper are given in [5].

*Example 3.* One could extend the previous pattern `SCAN_DRIVES` to account for alternative library calls, e.g. calling `FindFirstFileEx` instead of `FindFirstFile`:  
`SCAN_DRIVES = GetLogicalDriveStrings.GetDriveType.(FindFirstFile +FindFirstFileEx)`. This corresponds to the following SRS:

$$\begin{aligned} \text{FindFirstFile} &\rightarrow A \\ \text{FindFirstFileEx} &\rightarrow A \\ \text{GetDriveType}.A &\rightarrow B \\ \text{GetLogicalDriveStrings}.B &\rightarrow \text{SCAN\_DRIVES}. \end{aligned}$$

### 4.3 Abstracting with respect to a Set of Regular Behavior Patterns

Finally, we generalize abstraction to a set of behavior patterns. Indeed, in practice, a suspicious behavior can be expressed as the combination of several behavior patterns  $\mathcal{B}_i$ , each of them abstracting into a distinct symbol  $\lambda_i$ .

Throughout this paper, we denote by  $\Gamma$  the set of symbols representing abstractions of our behavior patterns and we extend  $\Sigma'$  to  $\Sigma \cup \Gamma$ .

Let  $\mathcal{C} = \{\mathcal{B}_i \mid \mathcal{B}_i \subseteq \Sigma^*\}_{1 \leq i \leq n}$  be a finite set of regular behavior patterns respectively associated with distinct symbols  $\lambda_i \in \Gamma$ .

As a relation  $\rightarrow_{\mathcal{B}_i}$  is defined on  $(\Sigma' \cup V_i)^*$ , the relation  $\rightarrow_{\mathcal{C}} = \bigcup_{1 \leq i \leq n} \rightarrow_{\mathcal{B}_i}$  is defined on  $(\Sigma' \cup \bigcup_{1 \leq i \leq n} V_i)^*$ . We extend trace language normalization to a set of regular behavior patterns, and the trace language  $L$  is now normalized with  $\rightarrow_{\mathcal{C}}$  into the abstract trace language  $L \downarrow_{\mathcal{C}}$ .

**Theorem 2.** *Let  $\mathcal{C}$  be a finite set of regular behavior patterns. If  $L$  is regular, then so is  $L \downarrow_{\mathcal{C}}$ .*

#### 4.4 Projecting the Abstract Trace Language on $\Gamma$

Once a trace language has been normalized with respect to some set of behavior patterns, details that have not been processed by the normalization have to be pruned. A completely abstracted trace language is indeed expected to be defined on  $\Gamma$ , in order to be compared afterwards to reference abstract behaviors. This pruning operation is performed by projecting the abstract trace language on  $\Gamma$ .

**Definition 5 ( $\Gamma$ -abstract trace language).** *Let  $\mathcal{C}$  be a set of regular behavior patterns. Let  $L$  be a trace language for some machine  $\mathcal{M}$ .  $\widehat{L} = L \downarrow_{\mathcal{C}}|_{\Gamma}$  is called the  $\Gamma$ -abstract trace language of  $\mathcal{M}$  with respect to  $\mathcal{C}$ .*

Once the abstraction is complete, the simplified traces describe functionality-related actions and are consequently more robust and less complex. As they represent several different implementations of the same behavior, they allow to deal with functional polymorphism.

We can then compare the language of  $\Gamma$ -abstracted traces to the behavior of some known malware or to some generic malicious behavior defined on  $\Gamma$ .

Finally, the whole abstraction process could be repeated, as in Martignoni et al’s layered architecture [19]. A first layer would look up behavior patterns defined in terms of raw analysis data. A second layer would look up behavior patterns defined in terms of patterns from the first layer, and so on. However, this case is encompassed in our formalism. It suffices to define behavior patterns from the final layer directly in terms of the raw analysis data, by composition of the regular SRSs defining patterns of the different layers. The resulting patterns thereby remain regular on  $\Sigma$ , so our formalism can still be applied.

## 5 Trace Abstraction Using Finite State Automata

When considering a single trace, the associated trace language is trivially represented by an automaton. But when this trace language describes the set of traces  $Traces(\mathcal{M})$  of some machine  $\mathcal{M}$ , this set is in general undecidable or at least non-regular, so no automaton can precisely represent it. Nevertheless, one may build a regular approximation of this trace language, which is usually twofold:

- the trace language is over-approximated, when replacing for instance  $a^n \cdot b^n$  sequences (stemming from program loops for instance) by  $a^* \cdot b^*$  sequences or when coping with obfuscation or dynamic analysis shortcomings. The resulting automaton then contains spurious traces (false positives).  
Formally, if  $Traces(\mathcal{M})$  is over-approximated by  $L$ , then  $Traces(\mathcal{M}) \subseteq L$ .
- the trace language is under-approximated, when some hidden code path is not discovered or when some uninteresting path is discarded. The resulting automaton then misses some traces (false negatives).  
Formally, if  $Traces(\mathcal{M})$  is under-approximated by  $L$ , then  $L \subseteq Traces(\mathcal{M})$ .

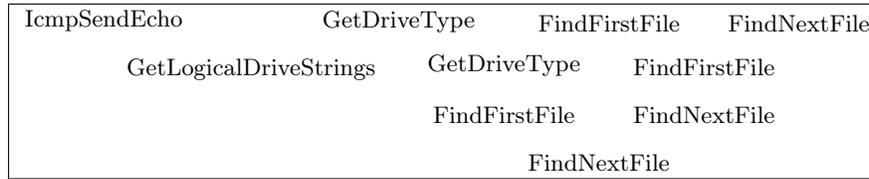
Thus a trace automaton represents a regular approximation of the trace language of some machine and can be defined as an automaton recognizing a part of the traces of  $\mathcal{M}$ .

**Definition 6 (Trace automaton).** *Given a machine  $\mathcal{M}$ , a trace automaton for  $\mathcal{M}$  with respect to  $\Sigma$  is a finite-state automaton  $A$  on  $\Sigma$  such that:*

$$\exists S \subseteq \Sigma^*, S \subseteq \text{Traces}(\mathcal{M}) \wedge S \subseteq \mathcal{L}(A).$$

In order to construct the trace automaton of a machine, one may either use a collection of captured traces or reconstruct the program machine code from its binary representation. In the first case, the automaton is built in such a way that the captured traces correspond to a path in the trace automaton. In the second case, the machine code of the program can be reconstructed using common techniques that combine static and dynamic analysis: it is then projected on the trace alphabet and the trace automaton is inferred from the code structure.

*Example 4.* Figure 1 shows a trace automaton for the Allapple.A excerpt representing the ping of a remote host and the scanning of local drives.



**Fig. 1.** Trace automaton for the Allapple.A excerpt

By Theorem 2, the abstraction problem for a regular trace language now amounts to computing an automaton recognizing the abstraction of this language. Construction of this automaton, which we call abstract trace automaton, is described in the proofs of the following theorems and uses a method proposed by Esparza et al. [13]. It consists in modifying the initial trace automaton by adding new transitions using the left hand sides of the rewriting rules and then intersecting it with an automaton recognizing the words in normal form with respect to our rewrite system.

Thus, the abstract trace automaton may be more complex than the initial one, as shown by the Allapple worm example. Abstraction of the trace automaton with respect to patterns `SCAN_DRIVES` and `PING`, where `PING = IcmpSendEcho` describes the ping of a remote host, gives the automaton of Figure 2.

**Theorem 3.** *Let  $A$  be a trace automaton and  $\mathcal{C} = \{\mathcal{B}_i\}_{1 \leq i \leq n}$  a set of behavior patterns recognized by regular SRSs  $\{R_{\mathcal{B}_i}\}_{1 \leq i \leq n}$ . Let  $|\mathcal{C}| = \sum_{1 \leq i \leq n} |R_{\mathcal{B}_i}|$ .*



**Fig. 2.** Abstract trace automaton for the Allapple.A excerpt

Then an automaton of size  $O(|A|)$  recognizing  $\mathcal{L}(A) \downarrow_{\mathcal{C}}|_{\Gamma}$  can be constructed in time  $O(|A|^3 \cdot |\mathcal{C}|^2)$  and space  $O(|A|^2 \cdot |\mathcal{C}|^2)$ .

The final abstraction of the Allapple.A excerpt, for  $\Gamma = \{\text{PING}, \text{SCAN\_DRIVES}\}$ , is depicted in Figure 3.



**Fig. 3.**  $\Gamma$ -abstract automaton for the Allapple.A excerpt

## 6 Application to Malware Detection

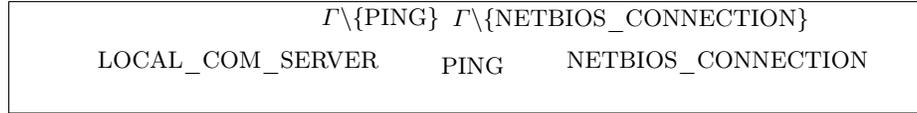
Using the abstraction framework defined in Section 4, malware detection now consists in computing the abstract trace language of some machine and comparing it to a database of malicious behaviors defined on  $\Gamma$ . These malicious behaviors either describe generic behaviors, e.g. sending spam or logging keystrokes, or behaviors of specific malware. According to our abstraction formalism, malicious behaviors are sets of particular combinations of behavior patterns abstractions.

**Definition 7.** A malicious behavior on  $\Gamma$ , or signature, is a language on  $\Gamma$ .

More specifically, a malicious behavior describes combinations of patterns, possibly interleaved with additional patterns which are irrelevant in these combinations. For instance, we define the signature for the Allapple worm as the following regular language, which explicitly allows interleaving of patterns that do not match the upcoming pattern:

$$\text{LOCAL\_COM\_SERVER} \cdot (\Gamma \setminus \{\text{PING}\})^* \cdot \text{PING} \cdot (\Gamma \setminus \{\text{NETBIOS\_CONNECTION}\})^* \cdot \text{NETBIOS\_CONNECTION}.$$

The automaton representing the signature of Allapple is given in Figure 4. Note that the `SCAN_DRIVES` pattern, which is present in the  $\Gamma$ -abstract trace automaton of the Allapple.A excerpt, does not appear here because the signature describes a common discriminating behavior exhibited by all samples of Allapple.



**Fig. 4.** Allapple signature

**Definition 8.** Let  $L_m$  be a malicious behavior on  $\Gamma$ . A machine  $\mathcal{M}$ , with a  $\Gamma$ -abstract trace language  $\widehat{L}$ , exhibits the malicious behavior  $L_m$  iff there exists  $v \in L_m$  and  $u \in \widehat{L}$  such that  $u = u_1vu_2$ , where  $u_1, u_2 \in \Gamma^*$ .

Thus,  $\mathcal{M}$  exhibits the behavior  $L_m$  if some subword of an abstract trace of  $\widehat{L}$  is in  $L_m$ . Our *malicious database* is then a set  $\mathcal{D}$  of malicious behaviors. A machine  $\mathcal{M}$  is *malicious* with respect to  $\mathcal{D}$  if it exhibits some malicious behavior of  $\mathcal{D}$ .

$\widehat{L}$  can be constructed either from a single captured trace or from a whole trace language. In the first case, the detection process is lighter. In the second case, we get a better description of the program behavior, so detection is more reliable.

When  $\widehat{L}$  is represented by an automaton  $A$  and  $L_m$  by an automaton  $A_m$ ,  $\mathcal{M}$  exhibits the behavior  $L_m$  when the following holds. Let  $A'_m$  be the automaton recognizing the set of words containing a subword in  $L_m$ :  $A'_m$  is constructed from  $A_m$  by adding loops labelled by symbols of  $\Gamma$  on the initial state and the final states of  $A_m$ . Then  $\mathcal{M}$  exhibits the behavior  $L_m$  iff  $\mathcal{L}(A) \cap \mathcal{L}(A'_m) \neq \emptyset$ .

The malicious database is then represented by an automaton  $A_{\mathcal{D}}$  which is the union of the automata  $A_m$  representing the signatures and may be minimized.

**Theorem 4.** Let  $\mathcal{D}$  be a set of regular malicious behaviors on  $\Gamma$ , recognized by an automaton  $A_{\mathcal{D}}$ . Let  $\mathcal{M}$  be a machine, with a  $\Gamma$ -abstract trace language recognized by an automaton  $A$ . Then deciding whether  $\mathcal{M}$  is malicious with respect to  $\mathcal{D}$  takes time:  $O(|A_{\mathcal{D}}|^2 \cdot |A|^2)$ .

Note that runtime detection could be efficiently implemented from these results, since all constructions ( $\Gamma$ -abstraction, automata intersection) are incremental.

Now, infection could be defined more intuitively and more generally in the following way: behavior patterns, instead of representing building blocks of a malicious behavior, could directly represent a malicious behavior, which would then be defined on  $\Sigma$ .

**Definition 9.** Let  $\mathcal{C}$  be a set of behavior patterns. Let  $\mathcal{M}$  be a machine, with a trace language  $L$ .  $\mathcal{M}$  is malicious if  $L \downarrow_{\mathcal{C}}|_{\Gamma} \neq \{\epsilon\}$ .

With the above detection model, we lose the expressive power of abstraction and the robustness of our detection model. But on the other hand, detection is performed by the sole process of normalization.

## 7 Implementation and Experiments

We implemented the presented techniques in a tool which is able to capture execution traces of a given program, to build its trace automaton, to abstract it with respect to a set of predefined behavior patterns and to compare it to a malware database.

*Setting up the detection environment* In order to avoid static analysis shortcomings and to ignore unreachable code, we use dynamic analysis to construct a trace automaton for some program loaded into a machine. The program is instrumented using Pin [3], which allows us to collect library calls along with their arguments while the program is running. Other instrumentation tools, like Dynamorio [1], could have been utilized with similar results. Instrumentation allows us to perform low-level analysis of the execution and to have a tight control over the program. In particular, our instrumentation tool handles threads, child processes and hooks, and features a simple data flow analyzer which relates call arguments to previously used data. In order to reduce the size of captured traces and to capture the behavior of a program at a source code level, we only collect library calls directly made from the program code, ignoring calls originating from libraries themselves.

When an execution trace is captured, we construct a trace automaton by associating a state with each different instruction pointer responsible for making a library call. Threads are handled by associating a state with a set of instruction pointers. Additional execution traces can be used to complete the trace automaton. Automata are manipulated with the OpenFST library [2].

Behavior patterns are defined after observing malicious execution traces and extracting basic sequences likely to be part of a malicious behavior. These patterns often define usual interactions of the program with the system or the network. Once extracted, a sequence either defines a new behavior pattern or extends an existing pattern.

The malware database is a collection of malicious behaviors, which is built from a set of generic signatures along with signatures of known malware, constructed using their  $\Gamma$ -abstract trace automata. The resulting database automaton is minimized in order to speed up detection.

*Experiments* To test our detection method, samples of malicious programs were collected using a honeypot<sup>1</sup> and identified using Kaspersky Antivirus.

We defined 40 behavior patterns, extended to allow data constraints expressing properties about the arguments of the calls (see [5]). These constraints are compatible with our formalism and amount to modify the trace automaton by

---

<sup>1</sup> The honeypot of the Loria's High Security Lab: <http://lhs.loria.fr>.

adding transitions denoting the verification of a constraint. Examples of such patterns include writing to system files, persisting oneself in the Windows registry or browsing drives and files. A complete abstraction example for the Agent.ah worm is also given in [5].

Three experimentation scenarios were defined. In a first setting, we define signatures for given malware families, by analysis and  $\Gamma$ -abstraction of a few samples. Samples from different variants of each family are then compared to the signatures: several of these variants are new, i.e. they were not considered when defining the signatures. This allows us to test the applicability of our approach and its robustness against mutation. In a second setting, a more general signature is defined for a common malicious behavior encountered in different malware families. Several malicious samples are then tested, in order to find out which samples exhibit this behavior. This allows us to test the expressive power of behavior abstraction. In a third setting, sane applications are tested, to ensure that the rate of false positives is low.

We tested the above settings on known malware families, among which Allaple, Virut, Agent, Rbot, Afcare and Mimail. In particular, our honeypot shows that Allaple, Virut and Agent are currently among the most active worms, which makes their analysis particularly relevant.

In the first scenario, we constructed a signature for Allaple (Figure 4) and abstracted samples from the variants a, b, d and e. Three of them were successfully matched to the signature. The fourth variant made use of Windows services, which our instrumentation tool does not currently handle. The same scenario was repeated for Virut, with its variants ai, ak, ao, n, q and r. Detection was again successful for four of them, although another technical shortcoming was encountered, stemming from the inability of our tool to instrument injected code called in remote threads.

For the second scenario, the signature is a generic malicious behavior describing malware replication, defined by the following patterns: injection of its own code in foreign live processes, duplication on logical drives, duplication on network shares, creating a local COM server, and registering an OLE component.

We then analyzed 39 variants from 16 malware families, among which: Agent, Mimail, Avron, Drefir... Samples of these variants were abstracted and compared to the previous behavior. Twenty of them were successfully matched.

For the last scenario, we abstracted traces of common applications, among which: Notepad, Firefox, VLC... All of them failed to exhibit a malicious behavior from the previous experiments.

## 8 Conclusion and Future Work

In this paper, we have presented a new malware detection approach using abstract representations of malicious behaviors to identify malicious programs. Programs to analyze, represented as trace languages, are abstracted by rewriting with respect to elementary behavior patterns, defined as regular string rewriting

systems. Abstractions are then compared to a database of abstract malicious behaviors, which describe combinations of the former patterns.

Abstraction is the key notion of our approach. Providing an abstracted form of program behaviors and signatures allows us to be independent of the program implementation, to handle similar behaviors in a generic way and thus to be robust with respect to existing and future variants. The strength of our technique lies also in the fact that abstract malicious behaviors are combinations of elementary patterns: this allows us to efficiently summarize and compact the possible combinations likely to compose suspicious behaviors. Moreover, malicious behaviors are easy to update since they are expressed in terms of basic blocks. Behavior patterns themselves, as they describe basic functionalities, are easier to update than if they had described more complex ones.

Behavior abstraction may also prove useful in similarity analysis of malicious programs. Like for detection, the use of abstract representations of programs makes the analysis resilient to functional polymorphism and to minor changes.

We plan to extend our approach in several directions. The first one is concerned with data flow analysis. The behavior patterns we define are unable to express that the actions they match are actually related. Data flow information would allow to address this question. While this issue is not very relevant when matching local behaviors (eg. writing data to a file), it becomes more important when specifying wider behaviors which spread over large execution sequences. Work is in progress to handle data flow by using model checking with temporal logic formulae with parameters.

The second extension is concerned with interleaved behavior patterns. When two behavior patterns are interleaved in an execution trace, our approach can only match one pattern out of the two since, while rewriting the first one, the second will be consumed. Although interleaving does not occur very often in practice since common behavior patterns are small and not very prone to interleaving, we intend to propose non-consuming approaches to pattern abstraction.

Also, we would like to improve the construction of a trace automaton approximating a trace language. When it is built from execution traces, meaningless paths are created by interference between relevant paths. This increases the matching possibilities between the abstract trace automaton and the malicious signatures, which impacts the precision of the detection. A solution would be to duplicate function calls in the automaton when appropriate.

Lastly, captured traces do not give an exhaustive view of all possible executions and some interesting behaviors may only be observed when some conditions are met. We could use existing tools identifying these conditions. Sage [14] and BitScope [8] use symbolic execution to systematically discover new execution paths. Moser et al. [20] also address this problem in the malware context by monitoring how certain inputs impact the control flow and by instrumenting these inputs to drive the execution.

## References

1. DynamORIO. <http://dynamorio.org>.

2. OpenFST. <http://www.openfst.org/>.
3. Pin. <http://www.pintool.org>.
4. Martin Apel, Christian Bockermann, and Michael Meier. Measuring similarity of malware behavior. In *IEEE Conference on Local Computer Networks*, pages 891–898. IEEE, October 2009.
5. Philippe Beaucamps, Isabelle Gnaedig, and Jean-Yves Marion. Behavior Abstraction in Malware Analysis - Extended Version. HAL-INRIA Open Archive Number inria-00509486.
6. J. Bergeron, M. Debbabi, J. Desharnais, MM. Erhioui, Y. Lavoie, and N. Tawbi. Static detection of malicious code in executable programs. In *Symposium on Requirements Engineering for Information Security*, 2001.
7. Guillaume Bonfante, Matthieu Kaczmarek, and Jean-Yves Marion. Architecture of a morphological malware detector. *Journal in Computer Virology*, 2008.
8. David Brumley, Cody Hartwig, Zhenkai Liang, James Newsome, Dawn Song, and Heng Yin. Automatically identifying trigger-based behavior in malware. *Botnet Detection*, 36:65–88, 2008.
9. Danilo Bruschi, Lorenzo Martignoni, and Mattia Monga. Detecting self-mutating malware using control-flow graph matching. In *Conference on Detection of Intrusions and Malware & Vulnerability Assessment*, volume 4064 of *Lecture Notes in Computer Science*, pages 129–143. Springer, 2006.
10. Mihai Christodorescu, Somesh Jha, Sanjit A. Seshia, Dawn Song, and Randal E. Bryant. Semantics-aware malware detection. In *IEEE Symposium on Security and Privacy*, pages 32–46. IEEE Computer Society, 2005.
11. Fred Cohen. Computer viruses: Theory and experiments. *Computers and Security*, 6(1):22–35, 1987.
12. Thomas Dullien and Rolf Rolles. Graph-based comparison of executable objects. In *Symposium sur la Sécurité des Technologies de l'Information et des Télécommunications*, 2005.
13. Javier Esparza, Peter Rossmanith, and Stefan Schwoon. A uniform framework for problems on context-free grammars. *Bulletin of the EATCS*, 72:169–177, 2000.
14. Patrice Godefroid, Michael Y. Levin, and David Molnar. Automated whitebox fuzz testing. In *Network Distributed Security Symposium*. Internet Society, 2008.
15. Carl A. Gunter. *Semantics of Programming Languages: Structures and Techniques*. MIT Press, 1992.
16. Grégoire Jacob, Hervé Debar, and Eric Filiol. Malware behavioral detection by attribute-automata using abstraction from platform and language. In *International Symposium on Recent Advances in Intrusion Detection*, volume 5758 of *Lecture Notes in Computer Science*, pages 81–100. Springer, 2009.
17. Johannes Kinder, Stefan Katzenbeisser, Christian Schallhart, and Helmut Veith. Detecting malicious code by model checking. In *Conference on Detection of Intrusions and Malware & Vulnerability Assessment*, volume 3548 of *Lecture Notes in Computer Science*, pages 174–187. Springer, 2005.
18. Baudouin Le Charlier, Abdelaziz Mounji, and Morton Swimmer. Dynamic detection and classification of computer viruses using general behaviour patterns. In *International Virus Bulletin Conference*, pages 1–22, 1995.
19. Lorenzo Martignoni, Elizabeth Stinson, Matt Fredrikson, Somesh Jha, and John C. Mitchell. A layered architecture for detecting malicious behaviors. In *International symposium on Recent Advances in Intrusion Detection*, volume 5230 of *Lecture Notes in Computer Science*, pages 78–97. Springer, 2008.

20. Andreas Moser, Christopher Kruegel, and Engin Kirda. Exploring multiple execution paths for malware analysis. In *IEEE Symposium on Security and Privacy*, pages 231–245. IEEE Computer Society, 2007.
21. R. Sekar, M. Bendre, D. Dhurjati, and P. Bollineni. A fast automaton-based method for detecting anomalous program behaviors. In *IEEE Symposium on Security and Privacy*, pages 144–155. IEEE Computer Society, 2001.
22. Prabhat K. Singh and Arun Lakhotia. Static verification of worm and virus behavior in binary executables using model checking. In *Information Assurance Workshop*, pages 298–300. IEEE Press, 2003.