

# *Semantic Trace-based Malware Variants Detection*

Khalid Alzarooni

CREST - DCS - UCL

April 6, 2011

# *Outline*

*Overview*

*Trace-based approach*

*Experiments*

# Overview

## *Malware Variants*

- Speed of evolution of malware partly driven by automatic generation of program variants
- Semantic equivalence tables used in malware, e.g. polymorphic and metamorphic malware
- These alter “local behaviour” of programs but larger scale behaviour is unchanged

## Malware Problem

Anoirel S. Issa

Symantec, UK (EICAR 2009)

*“Poly or metamorphic engines have some essential components that help them build highly obfuscated code. A single engine is able to produce unique variants that can reach millions.”*

Malware evolution:  $M_0 \rightarrow M_1 \rightarrow M_2 \rightarrow M_3 \rightarrow \dots$

Syntactic view:  $code_0 \not\approx code_1 \not\approx code_2 \not\approx code_3 \not\approx \dots$

## *Some Code Obfuscation Schemes*

[Beaucamps, 2007, Ször, 2005]

Label	Category	Obfuscation
gi	Garbage insertion	$\{\} \rightarrow \{C\}$
op	Opaque predicate	$\{\} \rightarrow \{P^{T/F}\}$
ec	Equivalent command	$\{op\} \rightarrow \{\bar{op}\}$
rr	Register renaming	$\{Rx\} \rightarrow \{Ry\}$
cs	Command split	$\{C\} \rightarrow \{C_x, C_y\}$
cm	Command merging	$\{C_x, C_y\} \rightarrow \{C_{xy}\}$
cr	Command reorder	$\{(C_x, C_y)\} \rightarrow \{(C_y, C_x)\}$
..	...	...

*Example: a program  $P$  and its semantically equivalent variant  $P'$*

$P :$ <hr style="border: 0.5px solid black; margin: 5px 0;"/> <i>a</i> $R0 := n$ <i>b</i> $R1 := m$ <i>c</i> $R2 := R1$ <i>d</i> $R3 := R2 + R0$ <i>e</i> $R4 := R1 + k$ <i>f</i> $R5 := 1$ <hr style="border: 0.5px solid black; margin: 5px 0;"/>	$\longrightarrow$	$P' :$ <hr style="border: 0.5px solid black; margin: 5px 0;"/> <i>a'</i> $R0 := n$ <i>cr</i> <sub>1</sub> $JMP \ rr_1$ <i>gi</i> <sub>1</sub> $R22 := R22 + 1$ <i>op</i> <sub>1</sub> $P^T \ JMP \ cm$ <i>rr</i> <sub>1</sub> $R11 := m$ <i>gi</i> <sub>2</sub> $R22 := R22 + 1$ <i>cr</i> <sub>2</sub> $JMP \ op_1$ <i>cm</i> $R3 := R11 + R0$ <i>e'</i> <sub>1</sub> $R4 := k$ <i>e'</i> <sub>2</sub> $R4 := R4 + R11$ <i>rr</i> <sub>2</sub> $R15 := 1$ <hr style="border: 0.5px solid black; margin: 5px 0;"/>
--	-------------------	--

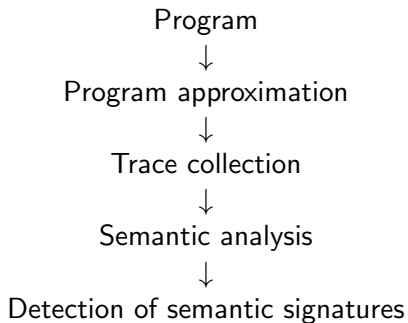
## Malware Problem

- To detect variants of a known malware
- Given two arbitrary programs is it possible to tell whether they are **semantically equivalent**?
- It is undecidable: not possible to devise an algorithm to produce “yes” or “no” detection answer [Cohen, 1987]

$$P' \stackrel{?}{\approx} P$$



## *Semantic trace-based*



## *Test scenarios*

### Results:

- Tested samples: Bho, Binom, Mobler, Telf, ...
- Most malware successfully matched, with  $k \geq 60\%$
- No false positives, similarity  $\leq 20\%$  (10 benign executables)
- 100% malware variants classification
- sig-w-slice: accuracy 30% and speed 26% in detection phase
- sig-wo-slice: 5:7 faster in sig. generation phase

# Trace-based approach

## Semantic trace-based

- Design a detector that can tell when two programs are *approximately* equivalent, which might often be good enough
- Approximate semantic equivalence is decidable
- Approximate a program's semantics  $\llbracket P \rrbracket$ 
  - CFG abstract traces (program paths) & test inputs
  - concrete & semantic traces

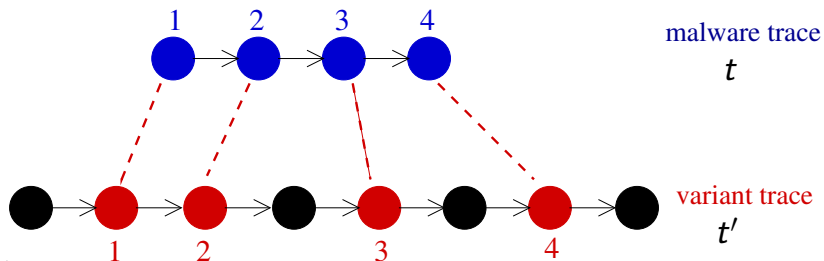
Malware evolution:  $M_0 \rightarrow M_1 \rightarrow M_2 \rightarrow M_3 \rightarrow \dots$

Syntactic view:  $code_0 \not\approx code_1 \not\approx code_2 \not\approx code_3 \not\approx \dots$

Semantic view:  $\llbracket M_0 \rrbracket \approx \llbracket M_1 \rrbracket \approx \llbracket M_2 \rrbracket \approx \llbracket M_3 \rrbracket \approx \dots$

## Semantic trace-based

- $M_1$  is a variant of  $M_0$  if  $\llbracket M_0 \rrbracket$  is **sub-sequence** of  $\llbracket M_1 \rrbracket$



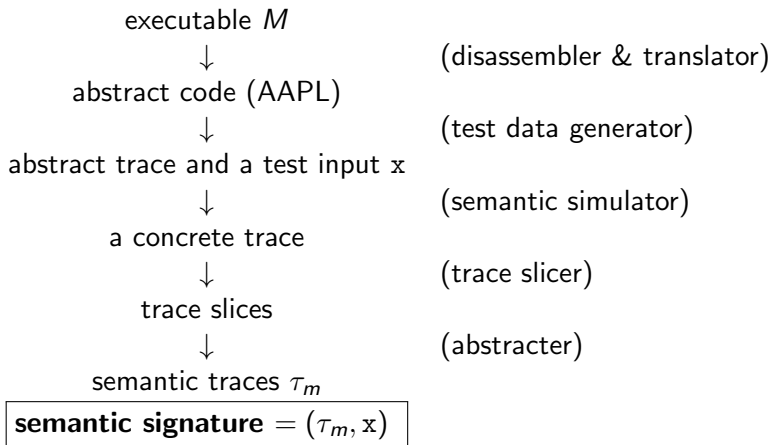
$$\forall t \in \llbracket M_0 \rrbracket, \exists t' \in \llbracket M_1 \rrbracket : t \prec t'$$

## *Semantic trace-based*

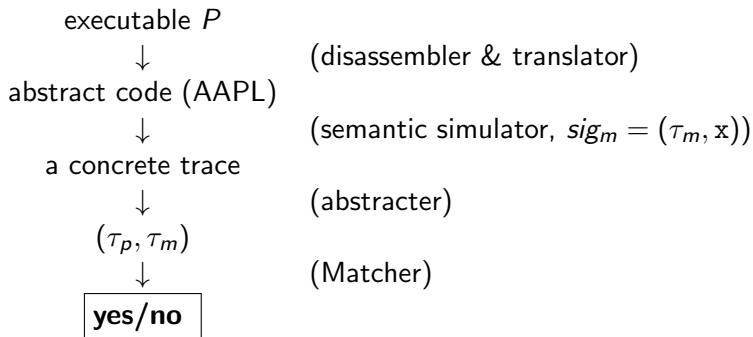
### Two phases:

1. Signature generation
2. Detection

## *Signature generation phase*



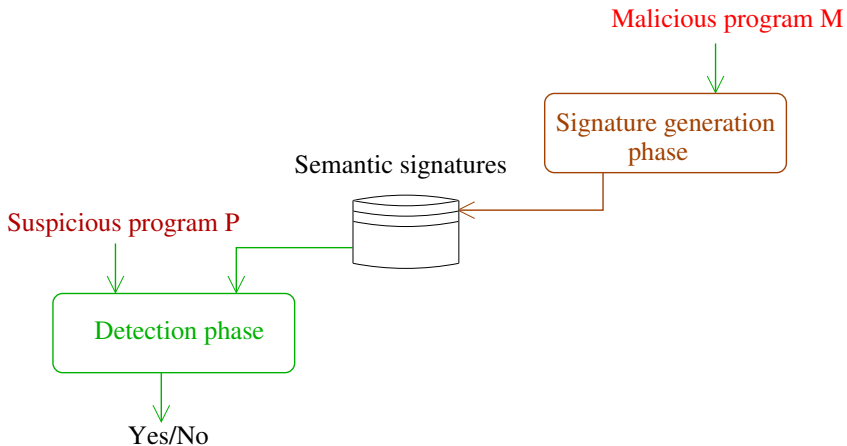
## Detection phase





# Experiments

## *Detector prototype*



## *Test scenarios*

We tested:

- Robustness against real in-the-wild variants
- Effectiveness of trace slicing in the signatures
- Fig. gen.& detection phases: sig-wo-slice vs. sig-w-slice
- False positives
- Classification of malware samples

## *Test scenarios*

### Results:

- Tested samples: Bho, Binom, Mobler, Telf, ...
- Most malware successfully matched, with  $k \geq 60\%$
- sig-w-slice: accuracy 30% and speed 26% in detection phase
- sig-wo-slice: 5:7 faster in sig. generation phase
- No false positives, similarity  $\leq 20\%$  (10 benign executables)
- 100% malware variants classification

## *Prototype limitation*

Technical shortcomings:

- Limited to viruses and worms
- Does not work for *dynamic* packed code and code with *anti-disassembly* techniques and
- Rely on tools to manually unpack (encrypted) and disassemble files

Thank you very much !



## References



Alzarouni, K., Clark, D., and Tratt, L. (2010).

Semantic malware detection.

Technical Report TR-10-03, Department of Computer Science,  
King's College London.



Beaucamps, P. (2007).

Advanced metamorphic techniques in computer viruses.

In *Proceedings of the International Conference on Computer,  
Electrical, and Systems Science, and Engineering - CESSE'07*.



Cohen, F. (1987).

Computer viruses: theory and experiments.

*Comput. Secur.*, 6(1):22–35.



Ször, P. (2005).

*The Art of Computer Virus Research and Defense*.

Addison-Wesley, Reading, Mass.

# Detector components



## Trace Semantics

- Trace semantics of a program is the set of all traces  $T$  that the program can produce
- A trace  $t \in T$  is a sequence of pairs of execution context  $\mathcal{X}$  and program syntax  $\mathcal{C}$
- Execution context: memory (locations) and environment (variables) values  $\mathcal{X} = \mathcal{E} \times \mathcal{M}$
- Program syntax: source code (commands)

$$\rho \in \mathcal{E} = \mathbf{R} \rightarrow \mathbb{Z}_{\perp} \quad (\text{environments})$$

$$m \in \mathcal{M} = \mathbb{Z} \rightarrow \mathbb{Z}_{\perp} \cup \mathbf{C} \quad (\text{memory})$$

$$\xi \in \mathcal{X} = \mathcal{E} \times \mathcal{M} \quad (\text{execution contexts})$$

$$S = \mathbf{C} \times \mathcal{X} \quad (\text{program states})$$

## *Trace Semantics*

- Signatures refer to exact program state
- Semantic signatures refer to values at particular memory locations and in registers that are observed to be constant across variants from the same malware family
- Detection: environment-memory traces of  $M$  that are contained (subtraces) of environment-memory traces of  $M'$

## Semantic Simulator

Not “live” testing

Evaluate abstract trace and collect concrete traces

Semantics of Actions:

$$\hat{\mathbf{A}} : \mathbf{A} \times \mathcal{X} \rightarrow \mathcal{X}$$

$$\hat{\mathbf{A}}[\mathbf{R} := \mathbf{E}]\xi = (\rho', m) \quad \text{where } \xi = (\rho, m) \text{ and } \rho' = \rho(\mathbf{R} \mapsto \hat{\mathbf{E}}[\mathbf{E}]\xi)$$

$$\hat{\mathbf{A}}[\mathbf{*R} := \mathbf{E}]\xi = (\rho, m') \quad \text{where } \xi = (\rho, m) \text{ and } m' = m(\rho(\mathbf{R}) \mapsto \hat{\mathbf{E}}[\mathbf{E}]\xi)$$

$$\hat{\mathbf{A}}[\mathbf{JMP} \ \mathbf{E}]\xi = (\rho', m) \quad \text{where } \xi = (\rho, m) \text{ and } \rho' = \rho(\mathbf{PC} \mapsto \hat{\mathbf{E}}[\mathbf{E}]\xi)$$

$$\hat{\mathbf{A}}[\mathbf{RTN}]\xi = (\rho', m) \quad \text{where } \xi = (\rho, m) \text{ and } \rho' = \rho(\mathbf{PC} \mapsto m(\rho(\mathbf{SP})), \mathbf{SP} \mapsto \mathbf{SP} + 1)$$

$$\hat{\mathbf{A}}[\mathbf{PUSH} \ \mathbf{E}]\xi = (\rho', m') \quad \text{where } \xi = (\rho, m) \text{ and } \rho' = \rho(\mathbf{SP} \mapsto \mathbf{SP} - 1) \text{ and} \\ m' = m(\rho(\mathbf{SP} - 1) \mapsto \hat{\mathbf{E}}[\mathbf{E}]\xi)$$

## Semantic Simulator

Not “live” testing

Evaluate abstract trace and collect concrete traces

Semantics of Commands:

$\hat{C} : S \rightarrow \Sigma(S)$  (determines transition relation between states)

$\hat{C}[[C_A]]\xi = (\xi', C')$  where  $\xi = (\rho, m)$ ,  $\xi' = \hat{A}[[A]]\xi$  and

$$C' = \begin{cases} m(\rho(PC)) & \text{if } A := \text{JMP} \cup \text{CALL} \cup \text{RTN} \\ m(\rho(PC + 1)) & \text{otherwise} \end{cases}$$

$\hat{C}[[C_B]]\xi = (\xi', C')$  where  $\xi = (\rho, m)$ , and

$$(\xi', C') = \begin{cases} \xi' = (\rho', m), \rho' = \rho(PC \mapsto \hat{E}[[E]]\xi), C' = m(\rho(\hat{E}[[E]]\xi)) & \text{if } \hat{B}[[B]]\xi = \text{true} \\ \xi' = \xi, C' = m(\rho(PC + 1)) & \text{otherwise} \end{cases}$$

## *TSA*lgo – Trace slicing

- $P \xrightarrow{\text{slice}} P'$  (semantically invariant subprogram wrt a criterion)
- $t \xrightarrow{\text{slice}} t'$  (semantically invariant subtrace wrt  $tsc$ )
- Trace slicing criterion  $tsc$ : recent definition points of variables in  $t$
- A conjecture: useful in the detection step for more accurate and efficient results.
- Effect is to shorten the trace and thus the signature

## *Signature matching*

$sig = (\tau_m, \mathbf{x})$  of  $M$  and  $\tau_p$  of  $P$ :

$$MD(sig, P) = \begin{cases} \text{yes} & \text{if } \tau_m \text{ is contained in } \tau_p \\ \text{no} & \text{otherwise} \end{cases}$$

Our assumption: some core semantic values in the two variants that would match with a high degree of similarity, indicating the likelihood of them being behaviourally the same.

## Signature matching

- we look for corresponding semantic traces of  $\tau_m$  in  $\tau_p$ ,
- a fuzzy matching to determine whether  $\tau_p$  corresponds semantically to  $\tau_m$

$$\text{semantics similarity measure} = \text{no. of mappings} / |\tau_m|$$

- We consider  $\tau_m$  is contained in  $\tau_p$  if the similarity measure is above a certain similarity threshold  $k$ ,

$$k \leq \text{similarity measure} \leq 100$$

$k$ : a large percentage of (desired) mappings