

A Semantics-based Approach to Malware Detection

Mila Dalla Preda – University of Verona, Italy

Mihai Christodorescu, Somesh Jha – University of Wisconsin, USA

Saumya Debray – University of Arizona, USA

17-19 Jan, POPL'07, Nice

A Few Basic Definitions

Malware represents malicious software.

Malware detector is a program \mathcal{D} that determines whether another program P is infected with a malware M .

$$\mathcal{D}(P, M) = \begin{cases} \text{True} & \text{if } \mathcal{D} \text{ determines that } P \text{ is infected with } M \\ \text{False} & \text{otherwise} \end{cases}$$

A Few Basic Definitions

Malware represents malicious software.

Malware detector is a program \mathcal{D} that determines whether another program P is infected with a malware M .

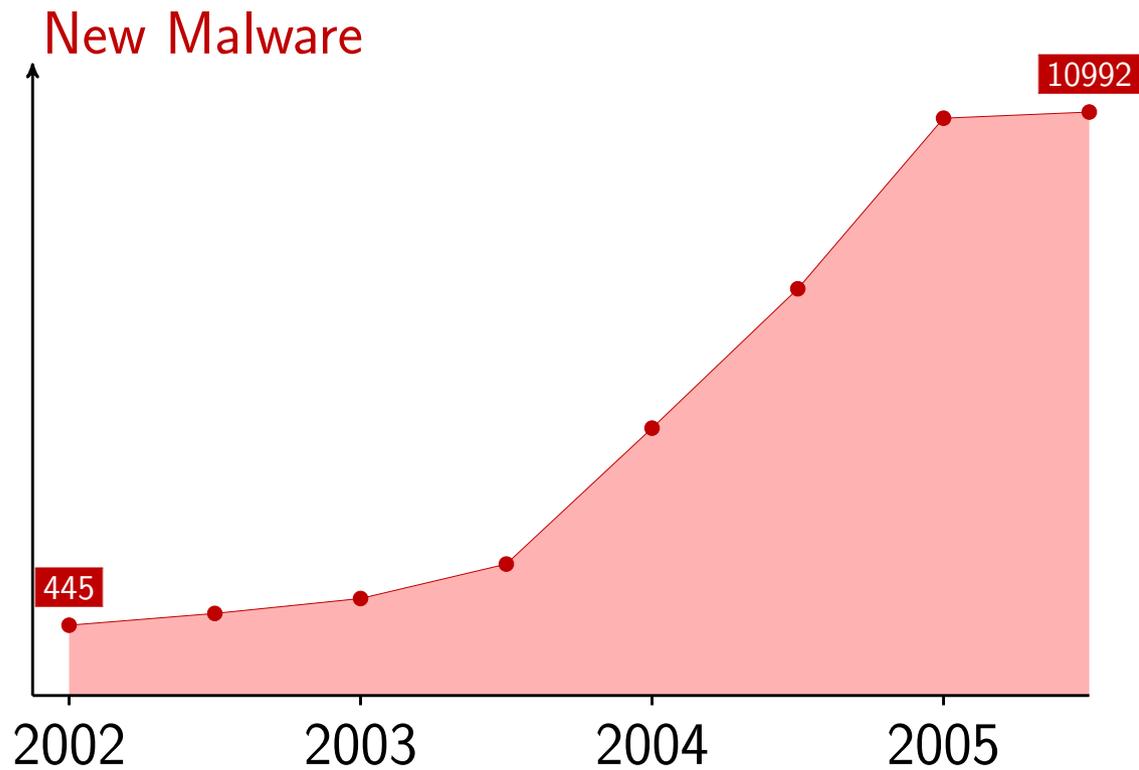
$$\mathcal{D}(P, M) = \begin{cases} \text{True} & \text{if } \mathcal{D} \text{ determines that } P \text{ is infected with } M \\ \text{False} & \text{otherwise} \end{cases}$$

An ideal malware detector detects all and only the programs infected with M , i.e., it is sound and complete.

- ⑥ Sound = no false positives (no false alarms)
- ⑥ Complete = no false negatives (no missed alarms)

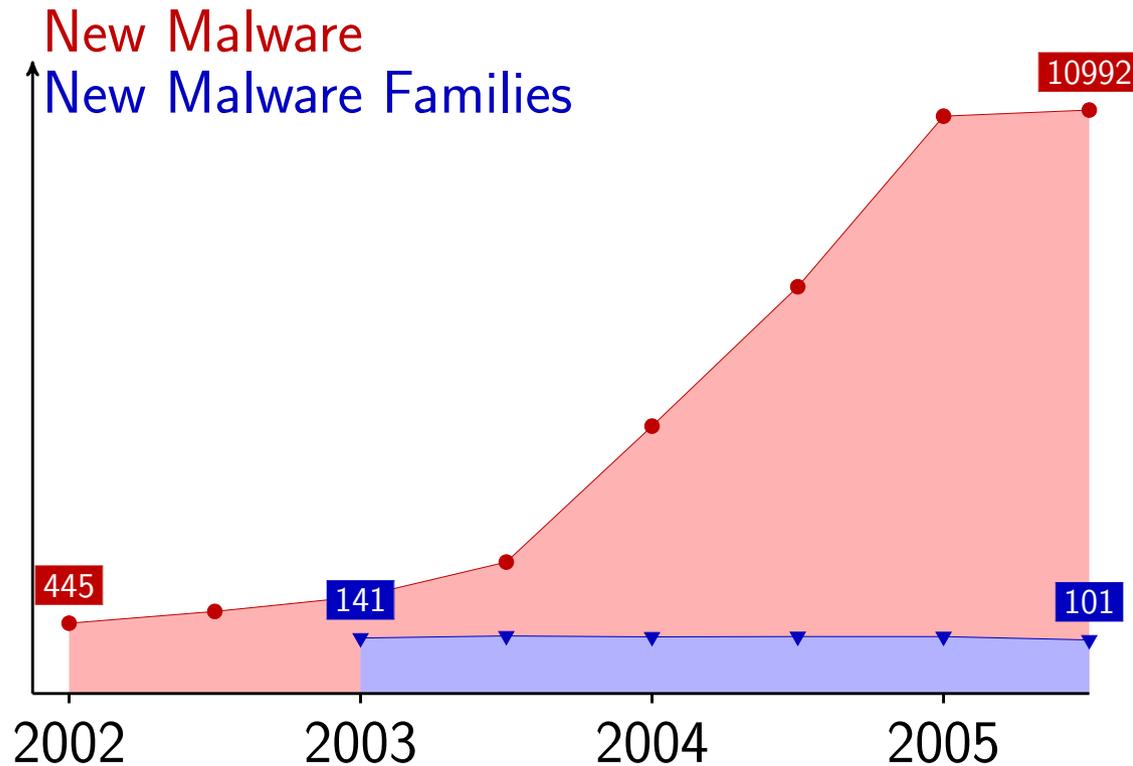
Malware Trends

There is more malware every year.



Malware Trends

There is more malware every year.



But the **number of malware families** has almost no variation.

Beagle family has 197 variants (as of Nov. 30).

Warezov family has 218 variants (as on Nov. 27).

The Malware Threat

Current detectors are **signature-based**:

P matches byte-signature **sig** \Rightarrow P is infected

Signature-based detectors, when sound, are not complete.

Malware writers use **obfuscation** to evade current detectors.

The Malware Threat

Current detectors are **signature-based**:

P matches byte-signature **sig** \Rightarrow P is infected

Signature-based detectors, when sound, are not complete.

Malware writers use **obfuscation** to evade current detectors.

Virus–antivirus “coevolution”

1. Malware writers create new, undetected malware.
2. Antimalware tools are updated to catch the new malware.
3. Repeat...

Common Obfuscations

- ⑥ Nop insertion
- ⑥ Register renaming
- ⑥ Junk insertion
- ⑥ Code reordering
- ⑥ Encryption
- ⑥ Reordering of independent statements
- ⑥ Reversing of branch conditions
- ⑥ Equivalent instruction substitution
- ⑥ Opaque predicate insertion
- ⑥ ... and many others...

Common Obfuscations

- ⑥ Nop insertion
- ⑥ Register renaming
- ⑥ **Junk insertion**
- ⑥ **Code reordering**
- ⑥ Encryption
- ⑥ Reordering of independent statements
- ⑥ Reversing of branch conditions
- ⑥ Equivalent instruction substitution
- ⑥ Opaque predicate insertion
- ⑥ ... and many others...

Obfuscation Example

(Pseudo-)Code:

```
mov eax, [edx+0Ch]
push ebx
push [eax]
call ReleaseLock
```

Obfuscation Example

(Pseudo-)Code:

```
mov eax, [edx+0Ch]
push ebx
push [eax]
call ReleaseLock
```

Obfuscated code (**junk**):

```
mov eax, [edx+0Ch]
inc eax
push ebx
dec eax
push [eax]
call ReleaseLock
```

Obfuscation Example

(Pseudo-)Code:

```
mov eax, [edx+0Ch]
push ebx
push [eax]
call ReleaseLock
```

Obfuscated code (**junk** + **reordering**):

```
mov eax, [edx+0Ch]
jmp +3
push ebx
dec eax
jmp +4
inc eax
jmp -3
call ReleaseLock
jmp +2
push [eax]
jmp -2
```

Solutions?

Recent developments based on **deep static analysis**:

- ⑥ Detecting Malicious Code by Model Checking [Kinder et al. 2005]
- ⑥ Semantics-Aware Malware Detection [Christodorescu et al. 2005]
- ⑥ Behavior-based Spyware Detection [Kirda et al. 2006]

Solutions?

Recent developments based on **deep static analysis**:

- ⑥ Detecting Malicious Code by Model Checking [Kinder et al. 2005]
- ⑥ Semantics-Aware Malware Detection [Christodorescu et al. 2005]
- ⑥ Behavior-based Spyware Detection [Kirda et al. 2006]

Lack of a formal framework for assessing these techniques.

Our Contributions

Challenges:

- ⑥ Many different obfuscations
- ⑥ Obfuscations are usually combined
- ⑥ Detection schemes usually rely on static/dynamic analyses

Our Contributions

Challenges:

- ⑥ Many different obfuscations
- ⑥ Obfuscations are usually combined
- ⑥ Detection schemes usually rely on static/dynamic analyses

A framework for assessing the resilience to obfuscation of malware detectors.

- ⑥ Obfuscation as transformation of trace semantics
- ⑥ Malware detection as abstract interpretation of trace semantics
- ⑥ Composing obfuscations vs. composing detectors

Two Worlds of Malware Detectors



Malware detector
on finite semantic structure

- ⑥ Disassembler
- ⑥ CFG construction
- ⑥ Other analyses

Two Worlds of Malware Detectors



Malware detector
on finite semantic structure

- ⑥ Disassembler
- ⑥ CFG construction
- ⑥ Other analyses



Malware detector
on trace semantics

Two Worlds of Malware Detectors



Malware detector
on finite semantic structure

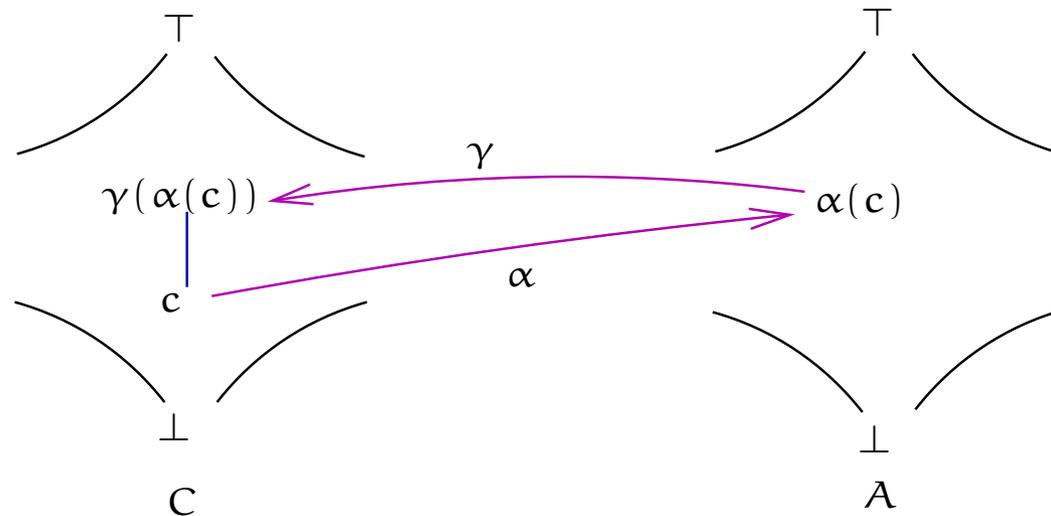
- ⑥ Disassembler
- ⑥ CFG construction
- ⑥ Other analyses



Malware detector
on trace semantics

Abstract Interpretation

Design approximate semantics of programs [Cousot & Cousot '77, '79].



Galois Connection: $\langle C, \alpha, \gamma, A \rangle$, A and C are complete lattices.

$\langle \text{Abs}(C), \sqsubseteq \rangle$ set of all possible abstract domains,

$A_1 \sqsubseteq A_2$ if A_1 is more concrete than A_2

Outline

- ⑥ Semantic Malware Detector
- ⑥ Soundness and Completeness
- ⑥ Classifying Obfuscations
- ⑥ Composing Obfuscations
- ⑥ Proving Soundness and Completeness

Semantic Malware Detector

A program P is infected by malware M , denoted $M \hookrightarrow P$
if (a part) of P execution is similar to that of M :

Semantic Malware Detector

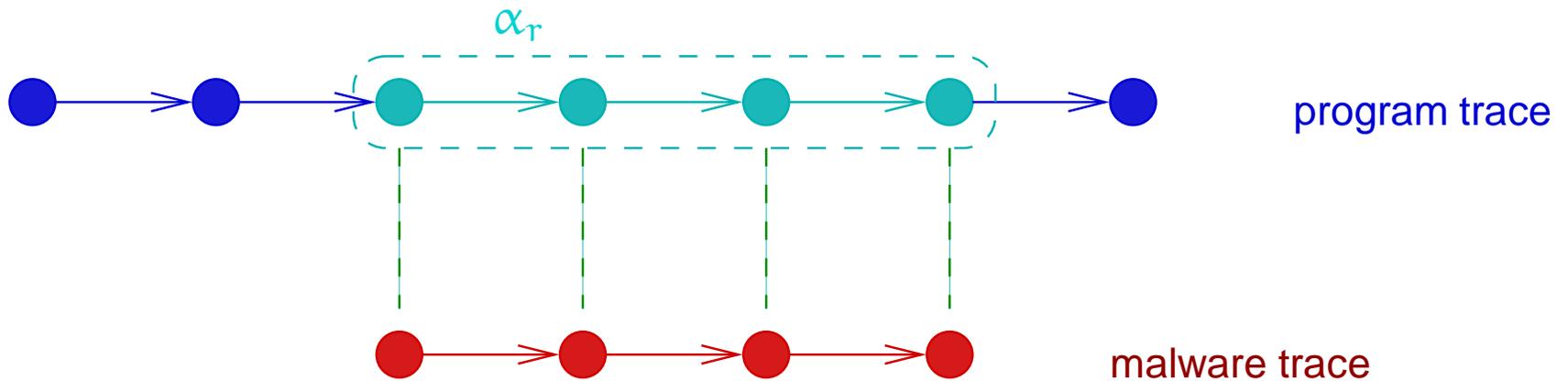
A program P is infected by malware M , denoted $M \hookrightarrow P$
if (a part) of P execution is similar to that of M :

$$S[M] \subseteq S[P]$$

Semantic Malware Detector

A program P is infected by malware M , denoted $M \hookrightarrow P$
if (a part) of P execution is similar to that of M :

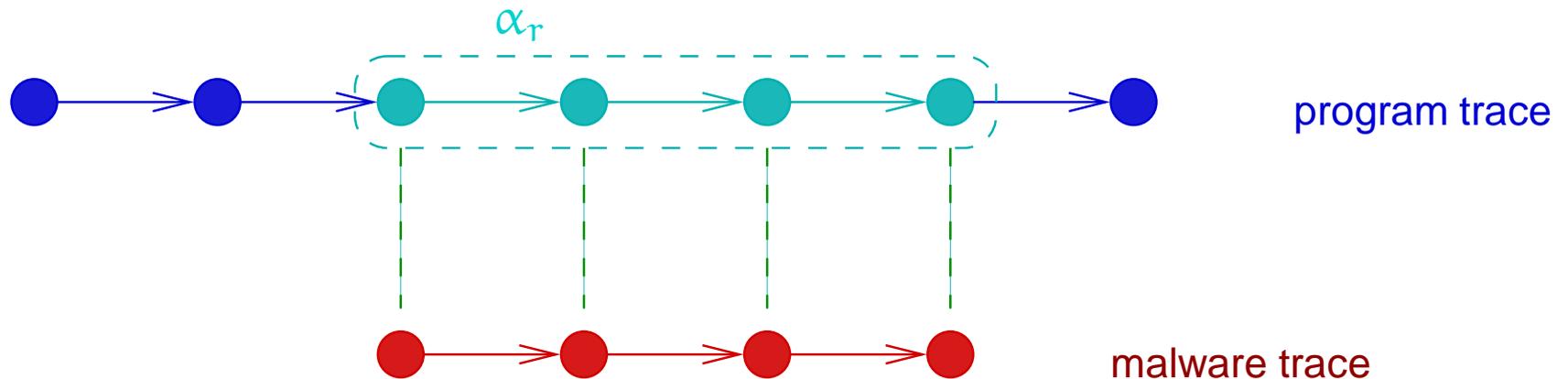
$$\exists \text{ restriction } r : S[[M]] \subseteq \alpha_r(S[[P]])$$



Semantic Malware Detector

A program P is infected by malware M , denoted $M \hookrightarrow P$
if (a part) of P execution is similar to that of M :

$$\exists \text{ restriction } r : S[[M]] \subseteq \alpha_r(S[[P]])$$



Vanilla Malware i.e. not obfuscated malware

Obfuscated Malware

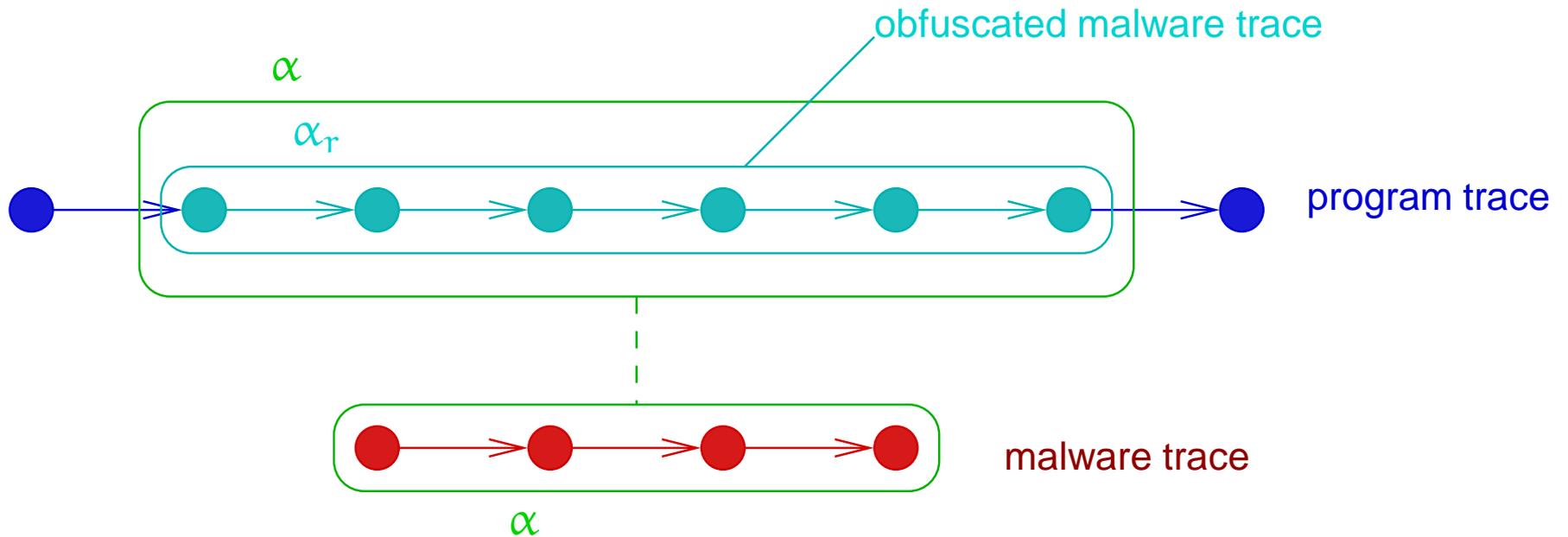
- ⑥ $\mathcal{O} : \mathbb{P} \rightarrow \mathbb{P}$ obfuscating transformation
- ⑥ $\alpha : Sem \rightarrow A$ abstraction that discards the details changed by the obfuscation while preserving maliciousness

$$\exists \text{ restriction } r : \alpha (S[[M]]) \subseteq \alpha (\alpha_r(S[[P]])$$

Obfuscated Malware

- ⑥ $\mathcal{O} : \mathbb{P} \rightarrow \mathbb{P}$ obfuscating transformation
- ⑥ $\alpha : Sem \rightarrow A$ abstraction that discards the details changed by the obfuscation while preserving maliciousness

\exists restriction $r : \alpha(S[M]) \subseteq \alpha(\alpha_r(S[P]))$



Sound vs. Complete

- ⑥ Precision of the Semantic Malware Detector (SMD) depends on α

Sound vs. Complete

- ⑥ Precision of the Semantic Malware Detector (SMD) depends on α
- ⑥ A SMD on α is **complete** w.r.t. a set \mathbb{O} of transformations if $\forall \mathcal{O} \in \mathbb{O}$:

$$\mathcal{O}(M) \hookrightarrow P \Rightarrow \begin{cases} \exists \text{ restriction } r : \\ \alpha(S[[M]]) \subseteq \alpha(\alpha_r(S[[P]])) \end{cases}$$

always detects programs that are infected (no false negatives)

Sound vs. Complete

⑥ Precision of the Semantic Malware Detector (SMD) depends on α

⑥ A SMD on α is **complete** w.r.t. a set \mathcal{O} of transformations if $\forall \mathcal{O} \in \mathcal{O}$:

$$\mathcal{O}(M) \hookrightarrow P \Rightarrow \begin{cases} \exists \text{ restriction } r : \\ \alpha(S[[M]]) \subseteq \alpha(\alpha_r(S[[P]]) \end{cases}$$

always detects programs that are infected (no false negatives)

⑥ If α is **preserved** by \mathcal{O} then the SMD on α is **complete** w.r.t. \mathcal{O} .

Sound vs. Complete

⑥ Precision of the Semantic Malware Detector (SMD) depends on α

⑥ A SMD on α is **sound** w.r.t. a set \mathbb{O} of transformations if:

$$\left. \begin{array}{l} \exists \text{ restriction } r : \\ \alpha(S[M]) \subseteq \alpha(\alpha_r(S[P])) \end{array} \right\} \Rightarrow \exists \mathcal{O} \in \mathbb{O} : \mathcal{O}(M) \hookrightarrow P$$

never erroneously claims a program is infected (no false positives)

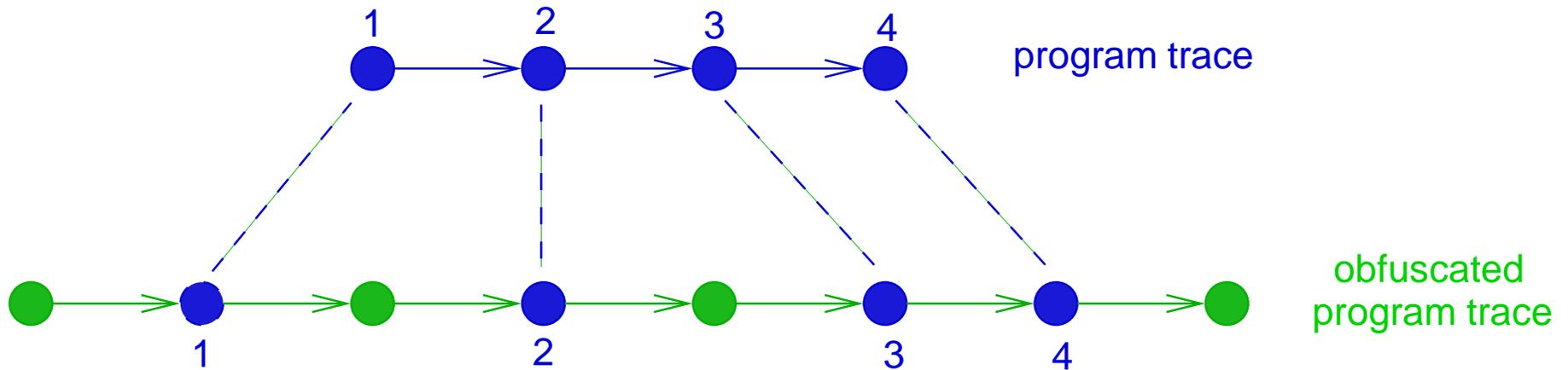
Outline

- ⑥ Semantic Malware Detector
- ⑥ Soundness and Completeness
- ⑥ Classifying Obfuscations
- ⑥ Composing Obfuscations
- ⑥ Proving Soundness and Completeness

Classifying Obfuscations

6 $\mathcal{O} : \mathbb{P} \rightarrow \mathbb{P}$ is a **conservative** obfuscation if

$\forall \text{trace1} \in S[\mathbb{P}], \exists \text{trace2} \in S[\mathcal{O}[\mathbb{P}}]$: trace1 is sub-sequence of trace2



Conservative Obfuscations

Abstraction α_c handles conservative obfuscations:

$$\alpha_c[X](Y) = X \cap \text{SubSequences}(Y)$$

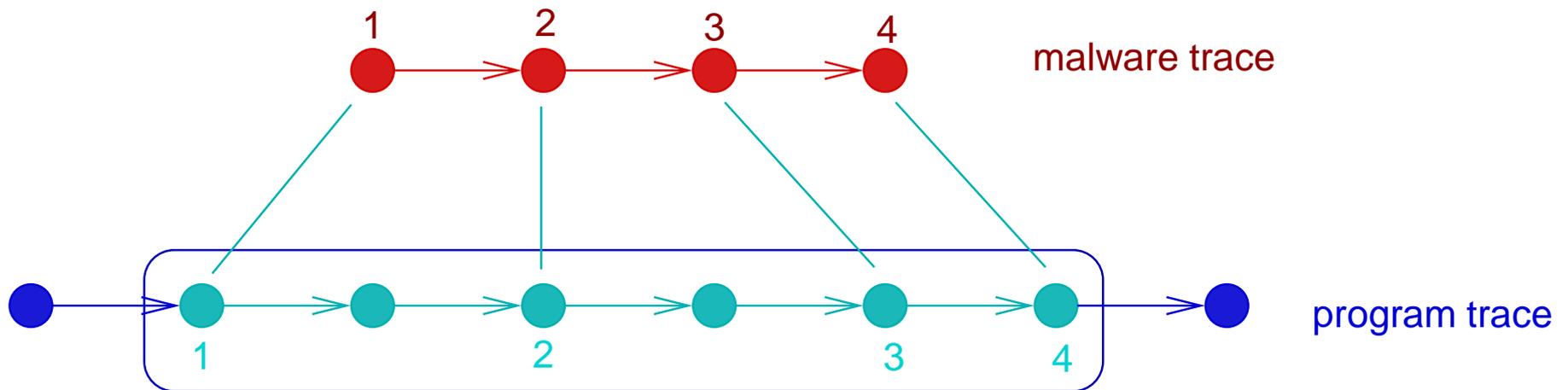
- 6 The SMD on α_c is **sound** and **complete** w.r.t. conservative obfuscations

Conservative Obfuscations

Abstraction α_c handles conservative obfuscations:

$$\alpha_c[X](Y) = X \cap \text{SubSequences}(Y)$$

- 6 The SMD on α_c is **sound** and **complete** w.r.t. conservative obfuscations

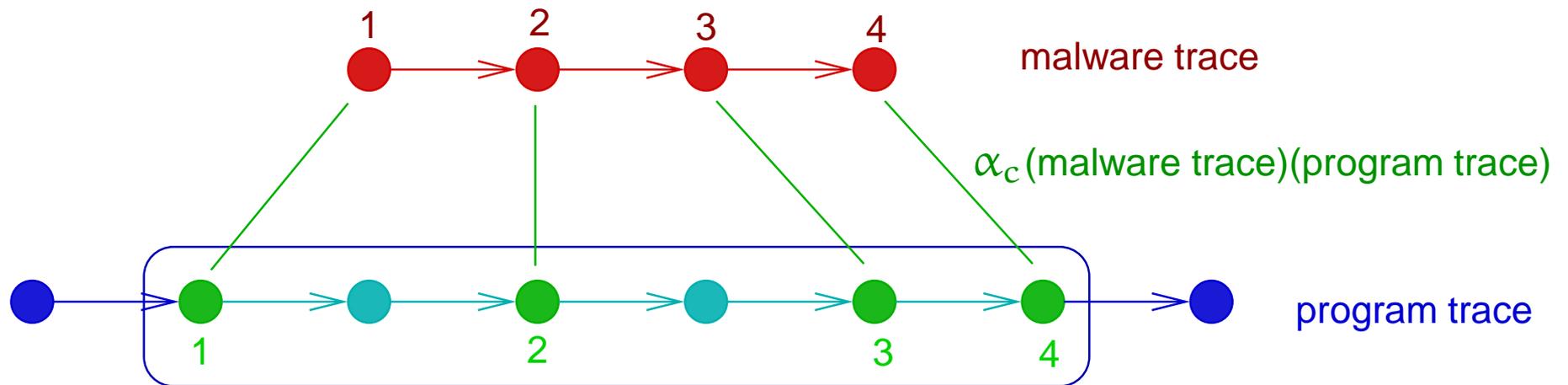


Conservative Obfuscations

Abstraction α_c handles conservative obfuscations:

$$\alpha_c[X](Y) = X \cap \text{SubSequences}(Y)$$

- 6 The SMD on α_c is **sound** and **complete** w.r.t. conservative obfuscations



Conservative Obfuscations

Abstraction α_c handles conservative obfuscations:

$$\alpha_c[X](Y) = X \cap \text{SubSequences}(Y)$$

- 6 The SMD on α_c is **sound** and **complete** w.r.t. conservative obfuscations

Abstraction α_c returns the set of malware traces that are subsequences of some program trace

Classifying Common Obfuscations

- ⑥ Nop insertion
- ⑥ Register renaming
- ⑥ Junk insertion
- ⑥ Code reordering
- ⑥ Encryption
- ⑥ Reordering of independent statements
- ⑥ Reversing of branch conditions
- ⑥ Equivalent instruction substitution
- ⑥ Opaque predicate insertion

Conservative Obfuscation Example

(Pseudo-)Code:

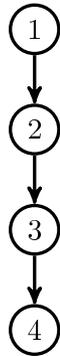
```
mov eax, [edx+0Ch]
push ebx
push [eax]
call ReleaseLock
```

Obfuscated code (junk + reordering):

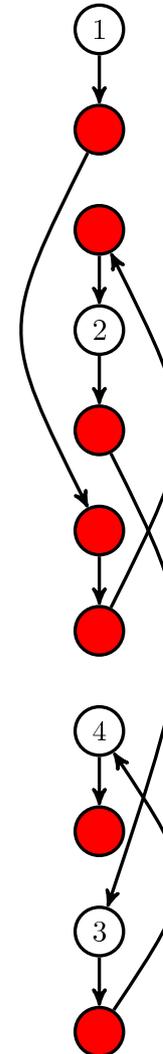
```
mov eax, [edx+0Ch]
jmp +3
push ebx
dec eax
jmp +4
inc eax
jmp -3
call ReleaseLock
jmp +2
push [eax]
jmp -2
```

Conservative Obfuscation Example

(Pseudo-)Code:

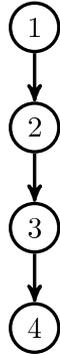


Obfuscated code (**junk** + **reordering**):



Conservative Obfuscation Example

(Pseudo-)Code:



Obfuscated code (**junk** + **reordering**):



Non-Conservative

Approach 1: Find a canonical transformation

Non-Conservative

Approach 1: Find a canonical transformation

(Pseudo-)Code:

```
mov eax, [edx+0Ch]
push ebx
push [eax]
call ReleaseLock
```

Obfuscated Code (Renaming):

```
mov edi, [eax+0Ch]
push ecx
push [edi]
call ReleaseLock
```

Non-Conservative

Approach 1: Find a canonical transformation

(Pseudo-)Code:

```
mov R1, [R2+0Ch]
push R3
push [R1]
call ReleaseLock
```

Obfuscated Code (Renaming):

```
mov R1, [R2+0Ch]
push R3
push [R1]
call ReleaseLock
```

Non-Conservative

- ⑥ Program infection: $M \hookrightarrow P$ if \exists restriction $r : S[[M]] \subseteq \alpha_r(S[[P]])$

Non-Conservative

⑥ Program infection: $M \hookrightarrow P$ if \exists restriction $r : S[[M]] \subseteq \alpha_r(S[[P]])$

Approach 2: Further abstractions

⑥ Interesting Malware States: $I \subseteq States[[M]]$:

$$M \hookrightarrow P \text{ if } \exists r : \alpha_I(S[[M]]) \subseteq \alpha_I(\alpha_r(S[[P]])$$



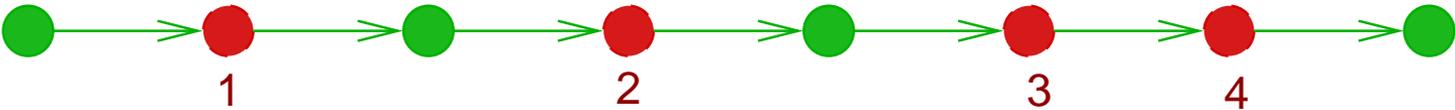
Non-Conservative

⑥ Program infection: $M \hookrightarrow P$ if \exists restriction $r : S[M] \subseteq \alpha_r(S[P])$

Approach 2: Further abstractions

⑥ Interesting Malware States: $I \subseteq States[M]$:

$$M \hookrightarrow P \text{ if } \exists r : \alpha_I(S[M]) \subseteq \alpha_I(\alpha_r(S[P]))$$



⑥ Interesting Malware Traces: $X \subseteq S[M]$

$$M \hookrightarrow P \text{ if } \exists r : X \subseteq \alpha_r(S[P])$$

Composition

- ⑥ Malware writers combine different obfuscations to avoid detection
- ⑥ The property of being conservative is preserved by composition
⇒ abstraction α_c
- ⑥ Under certain assumptions we can handle the composition of non-conservative obfuscations

Outline

- ⑥ Semantic Malware Detector
- ⑥ Soundness and Completeness
- ⑥ Classifying Obfuscations
- ⑥ Composing Obfuscations
- ⑥ Proving Soundness and Completeness

Proving Soundness/Completeness of MD

- ⑥ Identifying the class of obfuscators to which a malware detector is resilient can be a complex and error-prone task.
- ⑥ Obfuscators and detectors can be expressed on executions traces.

A detector is resilient to an obfuscator if it can “abstract away” the obfuscator’s effect on the program.

- ⑥ Case study: **Semantics-Aware Malware Detection Algorithm** proposed by [Christodorescu et al. 2005].
 - △ Complete for code reordering
 - △ Complete for junk insertion
 - △ Complete for variable renaming

Conclusions

- ⑥ Malware detection as abstraction of program semantics
vs.
Obfuscation as transformation of program semantics
- ⑥ We can now determine:
 - △ Whether a detector is resilient to a set of obfuscations
 - △ How complex a detector has to be to handle a given obfuscation
- ⑥ Open Problems:
 - △ Can we handle some interesting classes of non-conservative obfuscations?
 - △ How does one design a semantic detector based on trace semantics?
 - △ Connecting cryptographic and program analysis views of obfuscation

Thank you!