

# Context-Sensitive Analysis of Obfuscated x86 Executables

Arun Lakhotia(1), Davidson Boccardo(2), Anshuman Singh(1), and Aleardo Manacero Jr.(2)

(1)University of Louisiana at Lafayette, USA

(2)Paulista State University (UNESP), Brazil

PEPM 2010 (01/19/10)  
Madrid, Spain

# Disassembled binary with procedures: An example

Main:

```
L1:  PUSH 4
L2:  PUSH 2
L3:  CALL Max
L4:  PUSH 6
L5:  PUSH 4
L6:  CALL Max
L7:  PUSH 0
L8:  CALL ExitProcess
```

Max:

```
L9:  MOV  eax, [esp+4]
L10: MOV  ebx, [esp+8]
L11: CMP  eax, ebx
L12: JG   L14
L13: MOV  eax, ebx
L14: RET  8
```

# Context-sensitive interprocedural data-flow analysis - Classical methods

- **Call-string**
  - Sharir and Pnueli's  $k$ -call string method that maps a call string to its  $k$ -length suffix.
  - Emami *et al.*'s method of reducing recursive paths in a call string by a single node.
- **Procedure summary**
- **Inlining**

# Assumptions of call string based approaches

- The program uses special instructions like **call** and **ret** that can be identified and paired statically.
- Valid/invalid paths in ICFG can be described in terms of appropriate pairing of call-ret edges.

# Call and Ret are atomic

Call and Ret are **atomic** in the sense that they:

- Transfer control; and
- Change context

Call and Ret can be obfuscated using instructions that transfer control and change context separately. Call obfuscation can be employed by:

- **Malware writers**  $\Rightarrow$  to hide malicious behavior and to evade detection.
- **Software developers**  $\Rightarrow$  to protect intellectual property and to increase security.

# Call obfuscation using *push/ret* instructions

Main:

```
L1:  PUSH  4
L2:  PUSH  2
L3:  PUSH  offset [L6]
L4:  PUSH  offset [L13]
L5:  RET
L6:  PUSH  6
L7:  PUSH  4
L8:  PUSH  offset [L11]
L9:  PUSH  offset [L13]
L10: RET
L11: PUSH  0
L12: CALL  ExitProcess
```

Max:

```
L13: MOV   eax, [esp+4]
L14: MOV   ebx, [esp+8]
L15: CMP   eax, ebx
L16: JG     L18
L17: MOV   eax, ebx
L18: RET    8
```

# Call obfuscation using *push/jmp* instructions

Main:

```
L1:  PUSH  4
L2:  PUSH  2
L3:  PUSH  offset [L5]
L4:  JMP   Max
L5:  PUSH  6
L6:  PUSH  4
L7:  PUSH  offset [L9]
L8:  JMP   Max
L9:  PUSH  0
L10: CALL  ExitProcess
```

Max:

```
L11: MOV   eax, [esp+4]
L12: MOV   ebx, [esp+8]
L13: CMP   eax, ebx
L14: JG     L16
L15: MOV   eax, ebx
L16: RET    8
```



Classical call string based analyses are not directly applicable for context-sensitive analysis of binaries that have obfuscated calls. This is because:

- They are tied to semantics of procedure call and return statements of high-level languages, and therefore, call and ret instructions of assembly language.

**Objective:** Design of a context-sensitive analysis based on program semantics and abstract interpretation resilient from **call** and **ret** obfuscation attacks.

# Steps

- 1 Context abstractions (generic versions independent of ICFG based definitions)
- 2 Context-trace semantics (can not rely on ICFG based soundness results)
- 3 Language (a simple assembly language without call and ret)
- 4 Stack context (to model change of context)
- 5 Transfer of control (is modeled using value-set analysis)
- 6 Derive the context sensitive analyzer from context-insensitive one
- 7 Prove soundness of our analysis

# Generalized notion of contexts

- Opening and closing instructions are defined by:
  - $\langle \rangle \subseteq I$  - the set of instructions that open contexts.
  - $\rangle \rangle \subseteq I$  - the set of instructions that close contexts.
- For example, in the conventional interprocedural analysis, the set  $\langle \rangle$  contains the **call** instructions and  $\rangle \rangle$  contains the **ret** instructions.
- A **context-string** is a sequence of instructions that open contexts, represented by  $\langle \rangle^* \subseteq I^*$ .

- Let  $\mathbb{Q}^k$  represent the set of sequences of opening contexts of length  $\leq k$  and  $k + 1$  length sequences created by appending  $\top = \sqcup \mathbb{Q}$  to  $k$ -length sequences of opening contexts.
- An element of  $\mathbb{Q}^k$  is called a  $k$ -context. We can establish a map  $\alpha_k : \mathbb{Q}^* \rightarrow \mathbb{Q}^k$  as:

$$\alpha_k \nu \triangleq \begin{cases} \nu & \text{if } |\nu| \leq k \\ \nu_k.\top & \text{otherwise, where } \exists \nu' : \nu = \nu_k \wedge |\nu_k| = k. \end{cases}$$

- $\mathbb{Q}^*$  and  $\mathbb{Q}^k$  form a Galois insertion with the abstraction map  $\alpha_k$

- $\mathbb{Q}^\ell$  represent the set of sequence that open contexts with size  $\leq |\mathbb{Q}|$  and have cyclic sequence represented by  $+$ .
- For example, the term  $c^+$  represents all cyclic context strings from  $c$  to  $c$ .
- A map  $\alpha_\ell : \mathbb{Q}^* \rightarrow \mathbb{Q}^\ell$  can be defined such that  $\mathbb{Q}^*$  and  $\mathbb{Q}^\ell$  form a Galois insertion with the abstraction map  $\alpha_\ell$ .

# Examples of context abstractions

Context	2-Context	$\ell$ -Context
$C_2 C_1$	$C_2 C_1$	$C_2 C_1$
$C_2 C_3 C_2 C_1$	$C_2 C_3 \top$	$C_2^+ C_1$
$C_2 C_4 C_2 C_1$	$C_2 C_4 \top$	$C_2^+ C_1$
$C_2 C_4 C_2 C_3 C_2 C_1$	$C_2 C_4 \top$	$C_2^+ C_1$
$C_2 C_3 C_2 C_4 C_2 C_1$	$C_2 C_3 \top$	$C_2^+ C_1$
$C_3 C_2 C_4 C_2 C_1$	$C_3 C_2 \top$	$C_3 C_2^+ C_1$
$C_2 C_4 C_2 C_1$	$C_2 C_4 \top$	$C_2^+ C_1$
$C_5 C_2 C_4 C_2 C_1$	$C_5 C_2 \top$	$C_5 C_2^+ C_1$
$C_3 C_5 C_2 C_4 C_2 C_1$	$C_3 C_5 \top$	$C_3 C_5 C_2^+ C_1$
$C_5 C_5 C_2 C_4 C_2 C_1$	$C_5 C_5 \top$	$C_5^+ C_2^+ C_1$
$C_2 C_1$	$C_2 C_1$	$C_2 C_1$
$\epsilon$	$\epsilon$	$\epsilon$

- A **context-trace** is a pair of a context string and a trace  $(\nu, \sigma) \in (\mathbb{I}^* \times \Sigma^*)$ .
- The set of all context-traces of a program, denoted by  $\wp(\mathbb{I}^* \times \Sigma^*) \equiv \mathbb{I}^* \rightarrow \wp(\Sigma^*)$ , gives its **context-trace semantics**.



## Syntactic Categories:

$b \in \mathbf{B}$	(boolean expressions)
$e, e' \in \mathbf{E}$	(integer expressions)
$i \in \mathbf{I}$	(instructions)
$l, l' \in \mathbf{L} \subseteq \mathbb{Z}$	(labels)
$z \in \mathbb{Z}$	(integers)
$p \in \mathbf{P}$	(programs)
$r \in \mathbf{R}$	(references)

## Syntax:

$$e ::= l \mid z \mid r \mid *r \mid e_1 \text{ op } e_2$$
$$(op \in \{+, -, *, /, \dots\})$$
$$b ::= \text{true} \mid \text{false} \mid e_1 < e_2 \mid \neg b \mid$$
$$b1 \ \&\& \ b2$$
$$i ::= l : \text{esp} = \text{esp} + e \ . \ \text{eip} = e' \mid$$
$$l : \text{esp} = e \ . \ \text{eip} = e' \mid$$
$$l : * \text{esp} = e \ . \ \text{eip} = e' \mid$$
$$l : r = e \ . \ \text{eip} = e' \mid$$
$$l : *r = e \ . \ \text{eip} = e' \mid$$
$$l : \text{if } (b) \ \text{eip} = e; \ \text{eip} = l'$$
$$p ::= \text{seq}(i)$$

# Mapping Call and Ret in our language

- An instruction “*Call l*” may be mapped to the following sequence of instructions in our language:

$$l_0 : esp = esp - 1 \cdot eip = l_1$$

$$l_1 : *esp = l_2 \cdot eip = l$$

where  $l_2$  is the address of the instruction after the call instruction. It is not necessary that these two instructions appear contiguously in code.

- A *Ret* instruction may be mapped to the following instruction in our language:

$$l_0 : esp = esp + 1 \cdot eip = *esp$$

- **Idea:** To have the information about instructions that manipulate the stack pointer as a part of the context.
- The stack context can be described as the set of opening contexts and closing contexts represented by domains  $\langle \rangle_{asm} \subseteq I \times \mathbb{N}$  and  $\rangle \rangle_{asm} \subseteq I \times \mathbb{N}$  resp. that are defined as:

$$\langle \rangle_{asm} \triangleq \{(i, n) \mid \exists \delta, \delta' : \delta' \in (\mathcal{I} \ i \ \delta) \wedge (\delta' \ esp) = (\delta \ esp) - n\}$$

$$\rangle \rangle_{asm} \triangleq \{(i, n) \mid \exists \delta, \delta' : \delta' \in (\mathcal{I} \ i \ \delta) \wedge (\delta' \ esp) = (\delta \ esp) + n\}$$

- A context string is a sequence belonging to  $\langle \rangle_{asm}^*$ .  
Abstractions k-context and l-context can be applied to  $\langle \rangle_{asm}^*$  to reduce the complexity of the analysis.

# Transfer of control

- Upon execution of each instruction the instruction pointer register, *eip*, is updated with the label (a numerical value) of the next instruction to be executed.
- The value of the label may be computed from an expression involving values of registers and memory locations.
- We use Balakrishnan and Reps' **Value-Set Analysis (VSA)** to recover information about the contents of memory locations and registers. VSA uses the domain  $RIC = \mathbb{N} \times \mathbb{Z} \times \mathbb{Z}$  to abstract  $\wp(\mathbb{Z})$ .

# Derivation of a static analyzer

The analysis is derived from a chain of Galois connections linking the concrete domain  $\wp((I \times Store)^*)$  to the analysis domain  $I \rightarrow AbStore$ . The steps of the derivation are:

- The set  $\wp((I \times Store)^*)$ , called set of traces, is approximated to trace of sets, represented by  $(\wp(I \times Store))^*$ .
- The trace of sets is equivalent to  $(I \rightarrow \wp(Store))^*$ . This sequence of mapping of instructions to set of stores can be approximated to  $I \rightarrow \wp(Store)$ .
- Finally, a Galois connection between  $\wp(Store)$  and  $AbStore$  completes the analysis.

# Deriving the context-sensitive analyzer

Starting from concrete domain  $\llbracket^*_{asm} \xrightarrow{\Pi_{asm}} \wp(\Sigma^*)$  and the domain for Venable *et al.*'s context insensitive analyzer

$I \rightarrow R + L \rightarrow ASG \times RIC$ , we obtain our context sensitive analyzer analyzer  $\hat{\llbracket}^{\ell}_{asm} \rightarrow I \rightarrow R + L \rightarrow RIC$  using the following results:

- 1  $\llbracket^*_{asm} \sqsubseteq \hat{\llbracket}^{\ell}_{asm}$
- 2  $\wp(\mathbb{Z}) \sqsubseteq RIC$
- 3  $\llbracket^* \xrightarrow{\Pi} \wp(\Sigma^*) \equiv \wp(\Sigma^*)$

- The concrete context-trace semantics is given by the least fixpoint of the function

$$\mathcal{F}_c : \llbracket^*_{asm} \xrightarrow{\Pi_{asm}} \wp(\Sigma^*) \longrightarrow \llbracket^*_{asm} \xrightarrow{\Pi_{asm}} \wp(\Sigma^*), \text{ where } \Sigma = I \times R + L \rightarrow \mathbb{Z}.$$

- The context-trace semantics of the context-sensitive analyzer is given by the least fixpoint of the function  $\mathcal{F}^\# :$   
 $(\hat{\llbracket}^\ell_{asm} \rightarrow I \rightarrow R + L \rightarrow RIC) \longrightarrow (\hat{\llbracket}^\ell_{asm} \rightarrow I \rightarrow R + L \rightarrow RIC).$

## Lemma

$$\llbracket^*_{asm} \xrightarrow{\Pi_{asm}} \wp(\Sigma^*) \sqsubseteq \hat{\llbracket}^{\ell}_{asm} \rightarrow I \rightarrow R + L \rightarrow RIC.$$

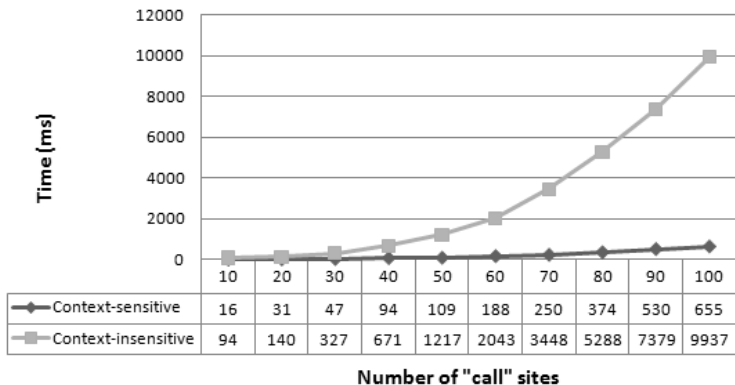
It follows from the lemma and the fixpoint transfer theorem that  $\mathcal{F}^{\#}$  is a sound approximation of  $\mathcal{F}_c$ .



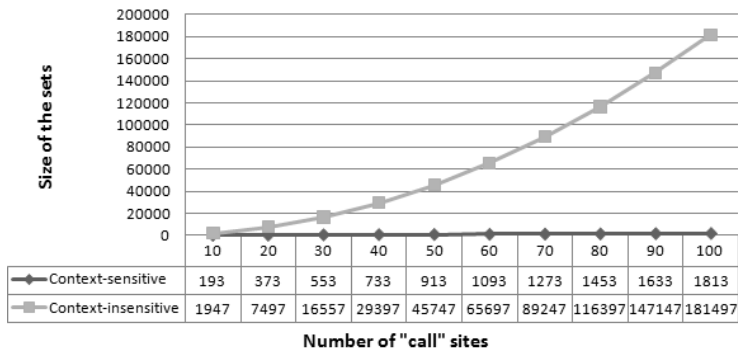
# DOC (Detector of Obfuscated Calls)

- We implemented our derived analysis in a tool called DOC.
- We studied the improvements in analysis of obfuscated code resulting from the use of our  $\ell$ -context-sensitive version of Venable *et al.*'s analysis against its context-insensitive version.
- We performed the analysis using two sets of programs:
  - Programs in the first set were **hand-crafted** with a certain known obfuscated calling structure.
  - The second set contains W32.Evol.a, a **metamorphic virus** that employs call obfuscation.

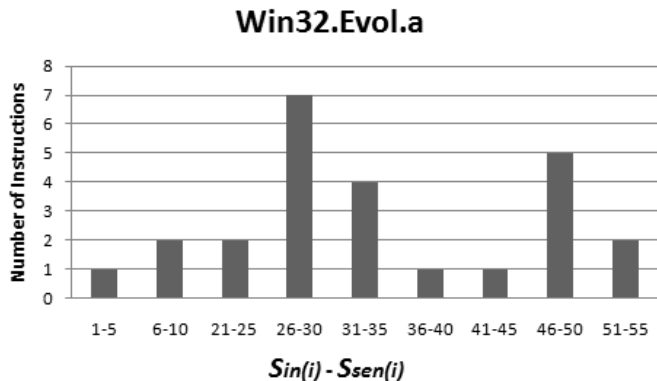
# Time evaluation



# Size of sets evaluation



# Histogram of evaluations for Win32.Evol.a



# Conclusions

- Developed a method for performing context sensitive analysis of binaries in which calling contexts cannot be discerned.
- Systematically derived generic versions of Sharir and Pnueli's k-suffix call-strings abstractions and Emami *et al.*'s strategy of abstracting calling-contexts (referred to as *l*-context in our work).
- Introduced the concept of stack-context, used in *lieu* of calling context, to perform context sensitive analysis of binaries that use call obfuscation.
- Proposed a general method for deriving sound context-sensitive analysis from context-insensitive one.