# Aversarial example againt SotA ML-based binary function classifiers

Gabriel Sauger

Loria

8th Franco-Japanese security workshop

8th Franco-Japanese security workshop

Gabriel Sauger (Loria) Aversarial example againt SotA ML-basec 1/23

#### Problem statement

Context:

- Target: Static binary **function** classification using machine-learning techniques.
  - Graph-based features
  - Assembly code analysis features (bytegrams, n-grams, strings...)

A classifier claim to extract the *semantics* of a given function, no matter the syntax (optimization / obfuscations).

Question: Can we build an obfuscator to automatically induce a misclassification in to SotA binary function classification models ?

# SotA papers in binary function classification

Marcelli et al (*How Machine Learning Is Solving the Binary Function Similarity Problem*) benchmarked several SotA binary function classifiers. We selected the best results:

- Asm2vec (assembly level 3-grams unsupervised)
- Google's GGSNN and GMN (graph and assembly level features)

All have near 0.90 on the benchmark performance metrics.

### Threat model

- Attacker can use their own compiler toolchain to compile the malicious binary.
- Attacker does not have access to C but knows the *probable* list of features used to construct C. This means that attacker knows the range of features on which the static analysis tools perform their analysis.
- The payload p is "smaller" than the target t
- Defender is restricted to static analysis based tools

#### Notations

- p: payload function
- t: target function, known to be benign to the defender
- *p<sub>asm</sub>*, *p<sub>src</sub>*, *p<sub>bin</sub>* is the *assembly code*, *source code*, *binary* version of p
- [[p]] is the *semantics* of p

**1** Run the CFG Merger:

$$(p_{src}, t_{src}) \rightarrow p'_{src}$$
 with  $\llbracket p \rrbracket = \llbracket p' \rrbracket$ 

8th Franco-Japanese security workshop

Gabriel Sauger (Loria) Aversarial example againt SotA ML-basec 7/23

Q Run the CFG Merger:

$$(p_{src}, t_{src}) 
ightarrow p_{src}'$$
 with  $\llbracket p 
rbracket = \llbracket p' 
rbracket$ 

8th Franco-Japanese security workshop

2 Compile to assembly code  $p'_{src}$  and  $t_{src}$ 

Run the CFG Merger:

$$(p_{src}, t_{src}) \rightarrow p'_{src}$$
 with  $\llbracket p \rrbracket = \llbracket p' \rrbracket$ 

- 2 Compile to assembly code  $p'_{src}$  and  $t_{src}$
- 8 Run the instruction frequency alignement algorithm

$$(p'_{asm}, t_{asm}) \rightarrow p''_{asm}$$
 with  $\llbracket p \rrbracket = \llbracket p'' \rrbracket$ 

Run the CFG Merger:

$$(p_{src}, t_{src}) \rightarrow p'_{src}$$
 with  $\llbracket p \rrbracket = \llbracket p' \rrbracket$ 

- 2 Compile to assembly code  $p'_{src}$  and  $t_{src}$
- Sun the instruction frequency alignement algorithm

$$(p'_{asm}, t_{asm}) \rightarrow p''_{asm}$$
 with  $\llbracket p \rrbracket = \llbracket p'' \rrbracket$ 

• Compile  $p''_{asm}$  to obtain  $p''_{bin}$  and feed it to the classifier.

CFG Merger: Principle

First step: CFG merge

8th Franco-Japanese security workshop

Gabriel Sauger (Loria) Aversarial example againt SotA ML-based 11/23

## CFG Merger: Principle

| Algorithm 1 Tree Matching algorithm, main func  | tion  |  |
|---|---|--|
| Require: (Prog, Prog)   |   |  |
| Ensure: Prog if match is possible, else None  |   |  |
| 1: function MERGE_PROG $(s, t)$   |   |  |
| 2: match $(s.cf, t.cf)$ with  |   |  |
| 3: $ (If, If) $   |   |  |
| 4: $res \leftarrow MERGE\_PROG(s.cf.body, t.cf.body)$   | (dy)  |  |
| 5: match res with   |   |  |
| 6: Prog   | ▷ We matched the bodies                                     |  |
| 7: $res2 \leftarrow \text{MERGE\_PROG}(s.cf.tail)$  | ,t.cf.tail)   |  |
| 8: match res2 with  | > And we matched the tails !                                |  |
| 9:   Frog   | ▷ And we matched the tans : a matched ::                    |  |
| 10: Teturn 1709(s.aata, re  | V Amazing, we have a match .)                               |  |
| 11: None  | Deport (a, t) > Here we need to incent a node from t in a   |  |
| 12: return MERGE_AFTER_INSERT(s,t) $\triangleright$ Here we need to insert a node from t in s |   |  |
| 13: end match   |   |  |
| 14: None  | (- t)   |  |
| 15: return MERGE_AFTER_INSERT(s,t)  |   |  |
| 16: end match   |   |  |
| 17:   (While, While)  | · Come Referention or more leader but with two WIL'         |  |
| 18: Same as previous case   | > Same disjonction as precedently but with two <i>w nue</i> |  |
| 19: $(If, While)$ or $(While, If)$ or $(End, If)$ or $(End, While)$                           |   |  |
| 20: return MERGE_AFTER_INSERT(s,t)  |   |  |
| 21: (End, End)  | <b>—</b> • • • • • • • • • • • • • • • • • • •              |  |
| 22: return Prog(s.data, End)  | ▷ Termination case, match :)                                |  |
| 23: (_, End)  |   |  |
| 24: return None   | $\triangleright$ Termination case, no match :(              |  |
| 25: end match   |   |  |
| 26: end function  |   |  |

Gabriel Sauger (Loria)

Aversarial example againt SotA ML-basec 12/23

<sup>2</sup> 8th Franco-Japanese security workshop

## CFG Merger: Principle

Second step: instruction matching

8th Franco-Japanese security workshop

Gabriel Sauger (Loria) Aversarial example againt SotA ML-basec 13/23

### Instruction-level obfuscation problem statement

We target an instruction-level feature based classifier C. Given the (x86) assembly code files of a payload p' and a target t, our goal is still to build a program p'' such that:

- $\llbracket p' \rrbracket = \llbracket p'' \rrbracket$  (p semantically equivalent to p')
- Idea: Aligning the histogram of instructions is enough to fool the classifier. (see Damasio et al *A Game-Based Framework to Compare Program Classifiers and Evaders*)

In this example, the score of  ${\cal C}$  is equal to 0 when the functions are dissimilar and 1 when they are the same semantically.

The selected features of an assembly instruction are:

- its operation (the mnemonic)
- its operands
  - register name
  - "immediate" if it's an address
  - reference if it's a "[eax + imm]"-like reference address

## Instructions obfuscation : hand-crafted examples

The steps are:

- Generate a bunch of  $\{(p_i, p'_i, t_i)\}_i$  with the merger.  $p \in sqlite$  and  $t \in curl$ .
- 2  $\forall i$ , build the **diff table i** of  $p'_i$  and  $t_i$ .
- Compute  $a2v(p_i'', t_i)$ ,  $a2v(p_i, t_i)$  and  $a2v(p_i'', p_i)$ . Compare.

Given a pair of function (p, t), we compute the **difference** between corresponding items in *Features*(p) and *Features*(t):

| feature             | Feat(p,t) = #t - #p |
|---------------------|---------------------|
| mox, rdx, rax       | 8                   |
| callimmediate       | 2                   |
|                     |                     |
| learax, [rbp + 0x3] | -5                  |

8th Franco-Japanese security workshop

Gabriel Sauger (Loria) Aversarial example againt SotA ML-based 17/23

# The Pair dataset

An eligible function is a function that:

- Has more than 5 basic blocks
- Is parsable by our CFG Merger

| Binary  | Nb compilable functions |
|---|-------------------------|
| source: <i>curl</i> - target: <i>openssl</i>  | 78                      |
| source: <i>curl</i> - target: <i>sqlite</i> 3 | 216                     |
| Total   | 294                     |

#### NOTE

*curl* and *openssl* binaries are included in the classifier's training dataset. *sqlite3* is not.

#### Benchmarks

Targets:

- GNN
- GMN
- asm2vec

<sup>O</sup> 8th Franco-Japanese security workshop

Gabriel Sauger (Loria) Aversarial example againt SotA ML-basec 19/23

# Instruction-level obfuscation: Misclassification problems

When we talk about misclassification, we refer to two problems:

#### "Query" problem

The defender has a database of function. We present him p'. The defender can then query his database, by looking at the top k highest-similarity functions according to C. Our goal then is to have t in the returned functions and not p.

One metric commonly used in this case is the MRR@k score.

#### "Triplet" problem

When given a triplet (p, t, p'), we want C to output that p' looks more similar to t than to p. Formally, if C(p, t) is low, our goal is:

$$\mathcal{C}(p',t) > \mathcal{C}(p',p) \text{ and } \mathcal{C}(p',t) > \mathcal{C}(p,t)$$

(1)

## Results: CFG Merger only

Results on the selected pairs, on the "Triplet problem":

| Dataset        | Classifier | "Good" triplets ratio |
|----------------|------------|-----------------------|
| (curl,sqlite3) | GNN        | 155/230 <b>(0.73)</b> |
| (curl,openssl) | GNN        | 63/78 <b>(0.81)</b>   |
| (curl,sqlite3) | GMN        | 38/211 <b>(0.18)</b>  |
| (curl,openssl) | GMN        | 16/78 <b>(0.21)</b>   |
| (curl,sqlite3) | a2v        | 6/211 <b>(0.03)</b>   |
| (curl,openssl) | a2v        | 1/78 <b>(0.01)</b>    |

# Results: CFG Merger + **random** instruction alignement algorithm

Results on the selected pairs, on the "Triplet problem":

| Dataset        | Classifier | "Good" triplets ratio |
|----------------|------------|-----------------------|
| (sqlite3,curl) | a2v        | (0.75)                |

While technically enough to have one misclassification with reasonnable "chance":

- Use semantic preserving instruction insertions
- Using save-regs / restore-regs functions
- "hiding" instructions in dead branches
- inserting sets of instructions equivalent to no-op
- compare to the Random Forest of Damasio et al.