

CSE 4800: Thesis

Binary Function Clone Detection Using MBA Simplification



Presented By -

Tasnimul Hasnat (190041113)

Sumit Alam Khan (190041207)

Md Abdullahil Kafi (190041224)

Supervised By -

Dr. Md Moniruzzaman,
Assistant Professor, IUT

Imtiaj Ahmed Chowdhury,
Lecturer, IUT

Table of Contents

1

Motivation

2

Relevant
Concepts

3

Existing
Literature

4

Problem
Statement

5

Proposed
Methodology

Motivation

- **Copyright infringement** is a serious issue in the software development industry.
- Verification and code **reuse identification** in closed-source applications is significantly difficult.
- Cybersecurity challenges increase as attackers **modify** the existing malwares to evade antivirus detection.
- Existing methodologies for finding function clone is mostly **Machine learning based** which is difficult to scale.
- A systematic method to find function clones from binary files could be **scalable** and easy to use.

Function Cloning

Function cloning is a method of **modifying** a function in such a way that the changes are limited to the internal code structure, **without altering** the function's overall **behavior**.

- Function cloning can be broadly categorized into two types:
 - Optimization
 - Obfuscation

Function Cloning (Optimization)

func1:

```
mov    -20[rbp], edi
mov    -24[rbp], esi
mov    -8[rbp], 10
mov    eax, -8[rbp]
xor    eax, -20[rbp]
sub    eax, -24[rbp]
mov    -4[rbp], eax
mov    eax, -24[rbp]
imul   eax, -4[rbp]
mov    edx, eax
mov    eax, -20[rbp]
add    eax, edx
ret
```

No Optimization

```
int func1(int a, int b){
    int x;
    int y;
    x = 10;
    y = (x ^ a) - b;
    return b * y + a;
}
```

Source
Function

func1:

```
mov    eax, edi
xor    eax, 10
sub    eax, esi
imul   esi, eax
lea    eax, [rsi+rdi]
ret
```

O1
Optimization

Function Cloning (Obfuscation)

```
int func1(int a, int b) {  
    int x;  
    int y;  
    x = 10;  
    y = ((x ^ a) & ~b) << 1)  
- ((x ^ a) ^ b);  
    return ((b * y - ~a) -1);  
}
```

Equivalent Obfuscated
Function

```
func1:  
    mov     edx, edi  
    xor     edx, 10  
    mov     eax, esi  
    not     eax  
    and     eax, edx  
    add     eax, eax  
    xor     edx, esi  
    sub     eax, edx  
    imul    esi, eax  
    not     edi  
    sub     esi, edi  
    lea     eax, -1[rsi]  
    ret
```

Optimized and
obfuscated

≈

```
func1:  
    mov     eax, edi  
    xor     eax, 10  
    sub     eax, esi  
    imul    esi, eax  
    lea     eax, [rsi+rdi]  
    ret
```

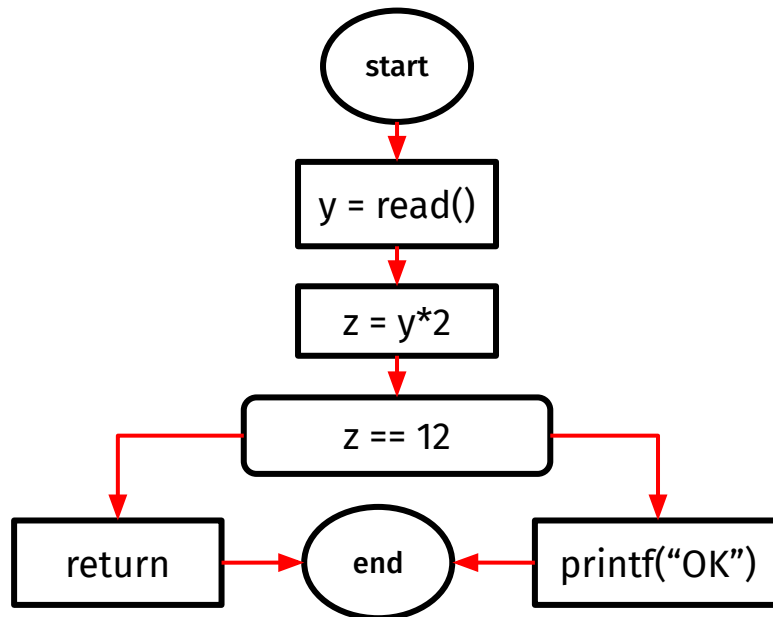
O1 Optimized
Function

Relevant Concepts

Control Flow Graph (CFG)

```
int f() {  
    ...  
    y = read();  
    z = y * 2;  
    if (z == 12) {  
        return;  
    } else {  
        printf("OK");  
    }  
}
```

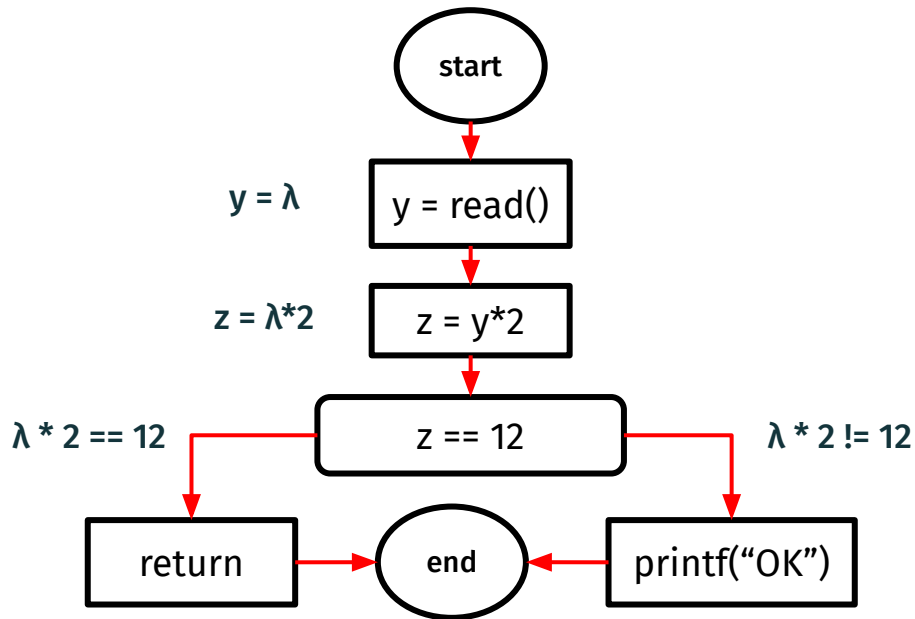
- Represents **flow of execution**
- **Nodes** → Basic block of codes
- **Edges** → Flow of control



Symbolic Execution

- Explores a **single** control flow **path** within a CFG.
- Execution is carried out using **symbolic inputs** rather than specific, concrete values.

```
int f() {  
    ...  
    y = read();  
    z = y * 2;  
    if (z == 12) {  
        return;  
    } else {  
        printf("OK");  
    }  
}
```



Mixed Boolean Arithmetic (MBA)

MBA Expression:

- ◆ Expression containing both traditional arithmetic operations (+, −, ×, ...) and bitwise boolean operators (∧, ∨, ¬, ⊕, ...)

$$\sum_{i \in I} a_i e_i(x_1, \dots, x_t)$$

where,

$x_t \rightarrow$ variables

$a_i \rightarrow$ constant coefficient

$e_i(x_1 \dots x_t) \rightarrow$ bitwise expressions

Figure: Formal Definition of MBA Expression¹

[1] **MBA-Blast: Unveiling and Simplifying Mixed Boolean-Arithmetic Obfuscation**
in the Proceedings of USENIX Security Symposium, 2021.

Mixed Boolean Arithmetic (MBA)

MBA Expression:

- ◆ A sequence of assembly instructions can be represented as MBA expressions.

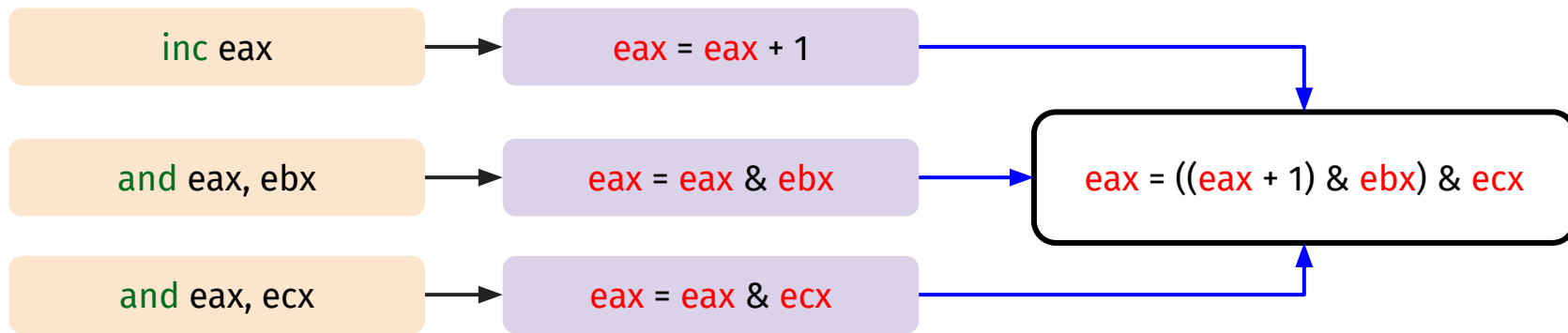


Figure: Assembly instructions to MBA expressions.

Existing Literature

Binary Function Clone Search in the Presence of Code Obfuscation and Optimization over Multi-CPU Architectures

ASIA CCS '23: Proceedings of the 2023 ACM Asia Conference on Computer and Communications Security, July 2023

Abdullah Qasem, Concordia University

Mourad Debbabi, Concordia University

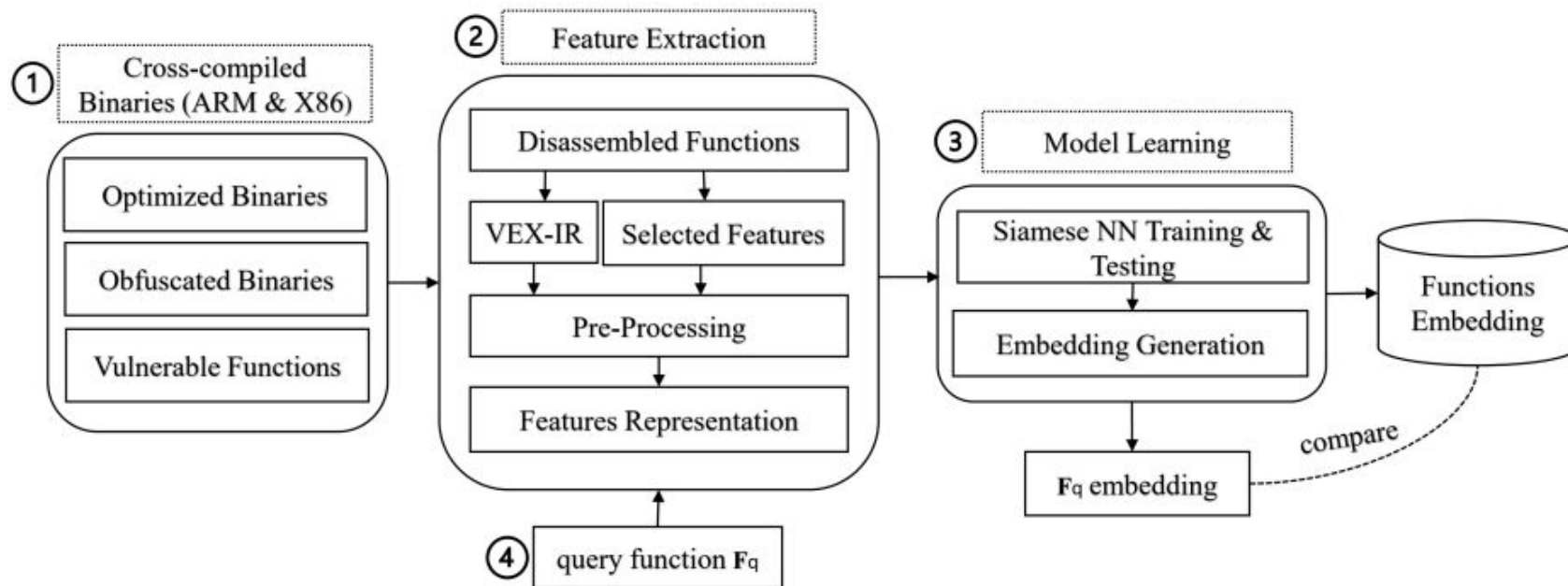
Bernard Lebel, Thales Research and Technologies

Marthe Kassouf, Hydro-Québec Research Institute

Contribution

- Developed and implemented **BinFinder**.
- Identified **resilient**, interpretable binary function features against code optimizations and obfuscations on **multi-CPU architectures**.
- Designed **Siamese** neural network architecture for training a model using proposed features to generate binary function **embeddings**.
- **Evaluated** BinFinder in 3 scenarios- x86 with different optimization and obfuscation, multi-cpu architecture with and without optimization and obfuscation.

BinFinder - Methodology



[2] **Binary Function Clone Search in the Presence of Code Obfuscation and Optimization over Multi-CPU Architectures**
ACM Asia Conference on Computer and Communications Security, July 2023

BinFinder - Evaluation

Approach	AUC				XM				
	XC	XC+XB	XA	XM	small	medium	large	MRR10	Recall@1
BinFinder	0.98	0.97	0.98	0.98	0.98	0.98	0.93	0.8	0.73
GMN_OPC-200_e16	0.86	0.85	0.86	0.86	0.89	0.82	0.79	0.53	0.45
GNN-s2v_GeminiNN_OPC-200_e5	0.78	0.81	0.82	0.81	0.84	0.77	0.79	0.36	0.28
SAFE_ASM-list_e5	0.8	0.8	0.81	0.81	0.83	0.77	0.77	0.17	0.27
Zeek	0.84	0.84	0.85	0.84	0.85	0.83	0.87	0.28	0.13
asm2vec	0.62	0.81	0.74	0.69	0.63	0.7	0.78	0.12	0.07

XA → Different Architecture, Same Obfuscation, Different Optimization.

XB → Same Architecture, Different Obfuscation.

XC → Different Architecture, Same Obfuscation, Same Optimization.

XM → Different Architecture, Different Obfuscation, Different Optimization.

- [2] **Binary Function Clone Search in the Presence of Code Obfuscation and Optimization over Multi-CPU Architectures**
ACM Asia Conference on Computer and Communications Security, July 2023

Qsynth- **A Program Synthesis based Approach for Binary Code** **Deobfuscation**

NDSS Symposium Binary Analysis Research Workshop, 2020

Robin David, *Quarkslab*

Luigi Coniglio, *University of Trento*

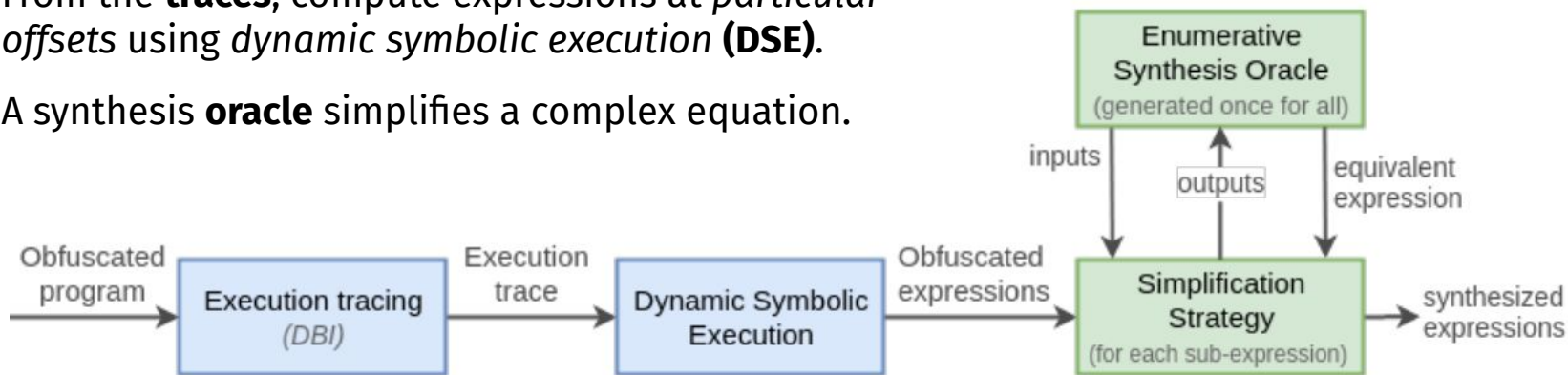
Mariano Ceccato, *University of Verona*

Contribution

- Introduced a combination of *dynamic symbolic execution*, *dataflow graph extraction*, and *program synthesis* against obfuscation.
- Proposed a **black-box synthesis** method based on *offline enumerative search*.
- Utilized precomputed lookup tables for near **constant-time** expression synthesis.
- **Iteratively** simplified complex obfuscated expressions.
- Empirical **validation** through a comparison with the deobfuscation tool **Syntia**, demonstrating superior accuracy and speed (20x faster).

Qsynth - Methodology

- 1 Uses *dynamic binary instrumentation (DBI)* to trace paths in the binary.
- 2 From the **traces**, compute expressions at *particular offsets* using *dynamic symbolic execution (DSE)*.
- 3 A synthesis **oracle** simplifies a complex equation.



[3] **Qsynth- A Program Synthesis based Approach for Binary Code Deobfuscation**
NDSS Symposium Binary Analysis Research Workshop, 2020

Qsynth - Evaluation

	Mean expr. size			Simplification			Mean scale factor			Sem.	Time			
	Orig	Obf _B	Synt	∅	Partial	Full	Obf _S /Orig	Synt/Obf _B	Synt/Orig		Sym.Ex	Synthesis	Total	per fun.
Syntia	/	/	/	52	0	448	/	/	/	/	/	/	34 min	4.08s
QSynth	3.97	203.19	3.71	0	500	500	x35.03	x0.02	x0.94	500	1m20s	15s	1m35s	0.19s

Orig, Obf_S, Obf_B, Synt are respectively original, obfuscated (source, binary level) and synthesized expressions

- [3] **Qsynth- A Program Synthesis based Approach for Binary Code Deobfuscation**
NDSS Symposium Binary Analysis Research Workshop, 2020

Simplifying Mixed Boolean-Arithmetic Obfuscation by Program Synthesis and Term Rewriting

**CCS '23: Proceedings of the 2023 ACM SIGSAC Conference on Computer and
Communications Security, November 2023. Pages 2351–2365**

Jaehyung Lee, Hanyang University, Ansan, Korea

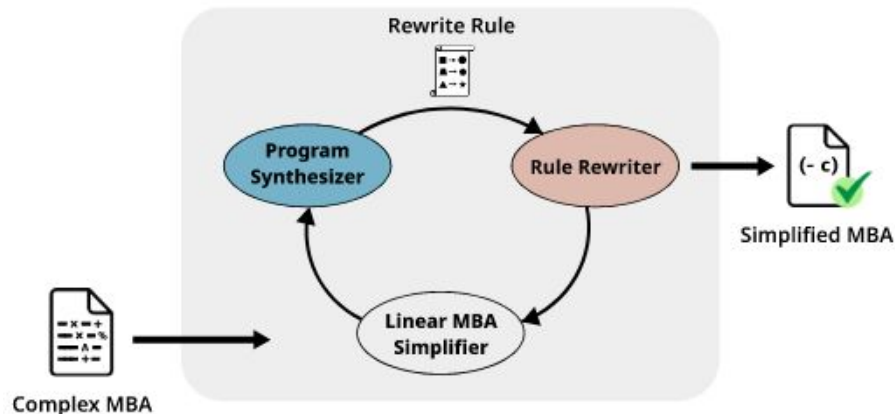
Woosuk Lee, Hanyang University, Ansan, Korea

Contributions

- Proposes a novel and versatile method called **ProMBA** for deobfuscating MBA expressions.
- Combines *program synthesis*, *term rewriting*, and *algebraic simplification* methods.
- Unlike existing techniques, ProMBA can deobfuscate a much **broader class** of MBA expressions.
- Implemented the proposed method in an **open source tool**.
- The experimental results show that ProMBA outperforms the state-of-the-art deobfuscation method, **MBASolver** by a large margin.

ProMBA - Methodology

- 1 Proposed method **first** simplifies **linear** MBA sub-expressions using an off-the-shelf deobfuscator.
- 2 **Recursively** simplifies **non-linear** sub-expressions by synthesizing simpler sub-expressions.
- 3 Applies the resulting rewrite rules to other sub-expressions **until no** further **simplification** is possible.

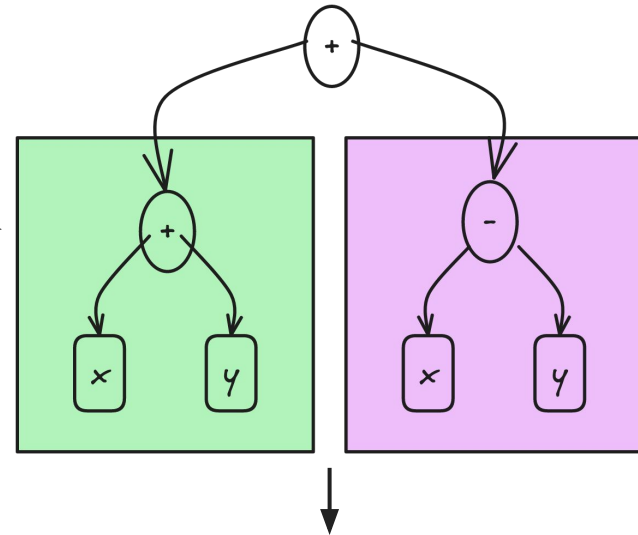
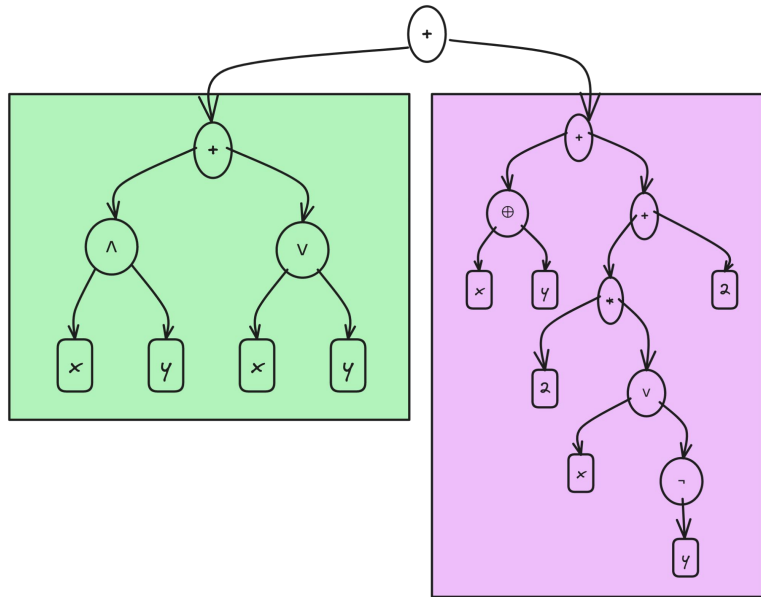


- [4] **Simplifying Mixed Boolean-Arithmetic Obfuscation by Program Synthesis and Term Rewriting**
ACM SIGSAC Conference on Computer and Communications Security, November 2023

MBA Simplification

Known Rules:

1. $x + y == (x \wedge y) + (x \vee y)$
2. $x - y == x \oplus y + 2 * (x \vee \neg y) + 2$



$$\{(x \wedge y) + (x \vee y)\} + \{x \oplus y + 2 * (x \vee \neg y) + 2\}$$

$$x + y + x - y \Rightarrow 2x$$

ProMBA - Evaluation

Dataset	Size (Avg.)			Success Rate			Time (Avg.)		
	PROMBA	MBASOLVER	SYNTIA	PROMBA	MBASOLVER	SYNTIA	PROMBA	MBASOLVER	SYNTIA
MBA-Solver	11.76	21.67	4.61	80.31%	25.16%	17.45%	65.64s	6.3s	12.95s
QSynth	17.48	77.71	4.72	62.8%	4.2%	22.8%	241s	64.83s	12.37s
Loki	3.51	866.25	3.1	97.2%	0.13%	74.4%	100.03s	347.77s	2.07s
Total	9.39	344.5	3.55	84.44%	13.19%	39.42%	100.36s	141.29s	8.81s

- [4] **Simplifying Mixed Boolean-Arithmetic Obfuscation by Program Synthesis and Term Rewriting**
ACM SIGSAC Conference on Computer and Communications Security, November 2023

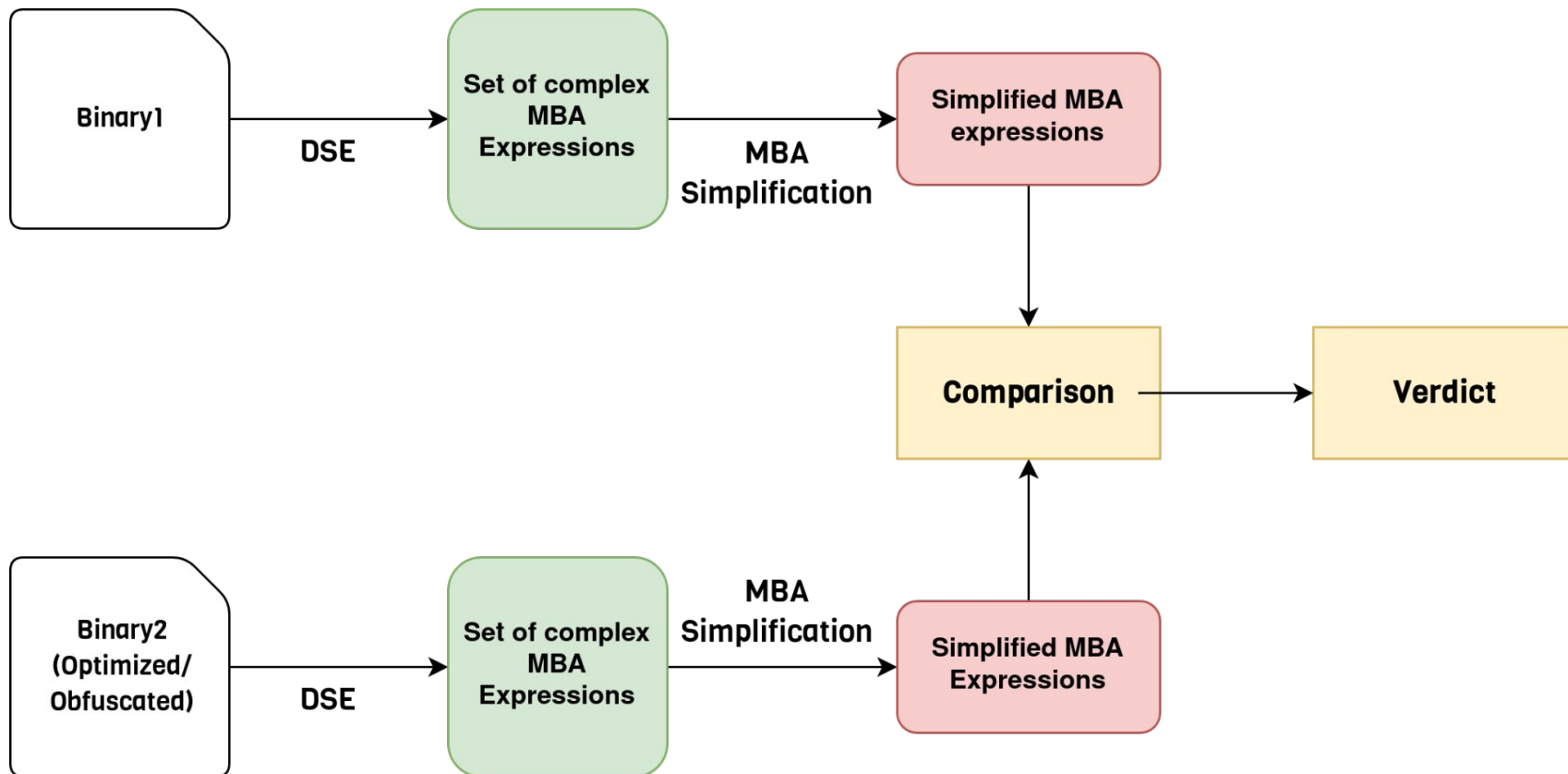
Limitations of existing literature

- Machine Learning based solutions incur unnecessary **overhead**.
- Existing synthesis methods spend time **deriving** the same expressions **again and again**.
- Works properly for only **small to medium** sized binaries.

Problem Statement

Given two binary files identify any function equivalency through extraction and simplification of MBA expressions

Proposed Methodology



Proposed Methodology

- 1 Create obfuscated binaries from **TREX** binary database.
- 2 Extraction of complex MBA expressions from obfuscated binary through dynamic symbolic execution.
- 3 Simplify the set of MBA expressions.
- 4 Check for **equivalency** between set of expressions using possible means e.g. SMT, fuzzing, program synthesis, graph isomorphism etc.

Challenges

- ◆ Extracting MBA expressions from branches and loops.
- ◆ Comparing equivalence of set of simplified expressions
- ◆ Choosing the appropriate form of verdict (binary vs similarity metrics).

Reference

- 1 **MBA-Blast: Unveiling and Simplifying Mixed Boolean-Arithmetic Obfuscation**
in the Proceedings of USENIX Security Symposium, 2021.
- 2 **Binary Function Clone Search in the Presence of Code Obfuscation and Optimization over Multi-CPU Architectures**
ACM Asia Conference on Computer and Communications Security, July 2023
- 3 **Qsynth- A Program Synthesis based Approach for Binary Code Deobfuscation**
NDSS Symposium Binary Analysis Research Workshop, 2020
- 4 **Simplifying Mixed Boolean-Arithmetic Obfuscation by Program Synthesis and Term Rewriting**
ACM SIGSAC Conference on Computer and Communications Security, November 2023

Thank You