

SIGMA: A Semantic Integrated Graph Matching Approach for Identifying Reused Functions in Binary Code[☆]

Saed Alrabaei, Paria Shirani, Lingyu Wang, Mourad Debbabi

National Cyber Forensics and Training Alliance Canada, Computer Security Laboratory, Concordia University, Montreal, Canada

Abstract

The capability of efficiently recognizing reused functions for binary code is critical to many digital forensics tasks, especially considering the fact that many modern malware typically contain a significant amount of functions borrowed from open source software packages. Such a capability will not only improve the efficiency of reverse engineering, but also reduce the odds of common libraries leading to false correlations between unrelated code bases. In this paper, we propose *SIGMA*, a technique for identifying reused functions in binary code by matching *traces* of a novel representation of binary code, namely, the *Semantic Integrated Graph (SIG)*. The *SIGs* enhance and merge several existing concepts from classic program analysis, including control flow graph, register flow graph, and function call graph into a joint data structure. Such a comprehensive representation allows us to capture different semantic descriptors of common functionalities in a unified manner as graph traces, which can be extracted from binaries and matched to identify reused functions, actions, or open source software packages. Experimental results show that our approach yields promising results. Furthermore, we demonstrate the effectiveness of our approach through a case study using two malware known to share common functionalities, namely, *Zeus* and *Citadel*.

Keywords:

Function Identification, Reverse Engineering, Binary Program Analysis, Malware Forensics, Digital Forensics

1. Introduction

The reverse engineering of binary code is generating significant interest among anti-virus companies, security experts, digital forensics consultants, law-enforcement agencies, national security agencies, etc. The objective of reverse engineering often involves understanding both the control and data-flow structures of the functions in the given binary code. However, this is usually a challenging task, because binary code inherently lacks structure due to the use of jumps and symbolic addresses, highly optimized control flow, varying registers and memory locations based on the processor and compiler, and the possibility of interruptions [2].

To assist reverse engineers in such a challenging task, automated tools for efficiently recognizing reused functions and their open source origins for binary code are highly desirable. This is especially true in the context of malware analysis, since modern malware are known to contain a significant amount of library code derived from either standard compiler libraries or open source software packages. The Flame malware, for instance, contains publicly available code packages, including SQLite and LUA [3]. Hence,

the ability to automatically identify reused functions may greatly enhance the effectiveness and efficiency of reverse engineering in such cases.

Existing techniques for identifying reused functions can be roughly categorized into static and dynamic approaches. In a static approach to function identification, different methods have focused on features at different levels (e.g., syntactical, semantical). For example, one existing technique counts mnemonics (opcode names, e.g., `add` or `mov`) in a sliding window over program text [13]. Another technique discovers exact and inexact clones in binaries through n-grams with normalization (linear naming of registers and memory locations) to address changes in names across different binaries [16]. Recently, an approach combines n-grams with small non-isomorphic sub-graphs of the control-flow graph to allow for structural matching [11]. More recently, another approach introduces tracelet-based code search in executables that attempts to statistically locate similar functions in the code base after translating the assembly instructions into an intermediate language [6]. While those techniques are not intended to address malware binaries, the authors in [15] identify shared software components to support malware forensics. In contrast to most static approaches that focus on one type of features, our approach combines different sources of information into one unified representation of binary code and thus has the potential of producing more accu-

[☆]This research is the result of a fruitful collaboration between the Computer Security Laboratory (CSL) of Concordia University, Defence Research and Development Canada (DRDC) Valcartier and Google under a DND/NSERC Research Partnership Program.

rate results. As to dynamic approaches, since they typically involve executing the code in order to detect the functionality, such approaches usually suffer from prohibitive runtime or exponential growth of execution paths [4, 9].

In this paper, we propose *SIGMA*, a technique for identifying reused functions in binary code by matching traces of a novel representation of binary code, namely, the semantic integrated graph (*SIG*). The *SIGs* enhance and merge several existing concepts from classic program analysis, including control flow graph, register flow graph, and function call graph, into a joint data structure. Such a comprehensive representation allows us to capture different semantic descriptors of common functionalities in a unified manner as traces of *SIG* graphs. Such *SIG* graph traces can then be extracted from binaries and matched, either exactly or approximately, to identify reused functions, actions, or open source software packages.

In summary, our contributions to the problem of identifying reused functions in binary code are as follows.

- We introduce the novel *SIG* representation of binary code to unify various semantic information, such as control flow, register manipulation, and function call into a joint data structure to facilitate more efficient graph matching.
- We define different types of traces such as normal traces, AND-traces, and OR-traces over *SIG* graphs, which are used for inexact matching. We carry out both exact and inexact matching between different binaries, where an exact matching applies to two *SIG* graphs with the same graph properties (e.g. number of nodes), whereas an inexact matching employs graph edit distance to measure the degree of similarity between two *SIG* graphs of different sizes.
- We evaluate our method by experimenting different variants of sort and encryption functions. Experimental results show that our method achieves similarity score close to an optimal similarity matching.
- Finally, we demonstrate the effectiveness of our approach through a case study using two known malware, which share common functionalities, namely, Zeus and Citadel.

The rest of the paper is organized as follows. Section 2 reviews several existing representations of binary code. Section 3 provides a detailed description of the main methodology. Section 4 evaluates the proposed approach and compares it to existing work. Section 5 describes our case study. Section 6 gives limitations and future directions. Section 7 reviews related work, and Section 8 draws conclusions.

2. Existing Representations of Binary Code

Numerous representations of binary code have been developed for different purposes of program analysis, such as

data flow analysis, control flow analysis, call graph analysis, structural flow analysis, register manipulation analysis, and program dependency analysis. While these representations have been designed primarily for analyzing binary code, they can certainly be employed to characterize the code. In particular, we focus on three representations that capture structural information, namely, control flow graph, register flow graph, and function call graph. These representations form the basis of our approach to identifying reused functions in binary code. For the sake of clarity, we introduce a running example to illustrate these representations using the following sample code (bubble sort).

```
void bubble_sort(int arr[], int size) {
    bool not_sorted = true;
    int j=0,tmp;
    while (not_sorted)
    {
        not_sorted = false;
        j++;
        for (int i = 0; i < size - j; i++){
            if (arr[i] > arr[i + 1]) {
                tmp = arr[i];
                arr[i] = arr[i + 1];
                arr[i + 1] = tmp;
                not_sorted = true;
            }
        }
        print_array(arr,5);
    }
}
```

2.1. Control Flow Graph

Control Flow Graphs (CFGs) have been used for a variety of applications, e.g., to detect variants of known malicious applications [5]. A CFG describes the order in which basic block statements are executed as well as the conditions that need to be met for a particular path of execution. To this end, basic blocks are represented by nodes connected by directed edges to indicate the transfer of control. It is necessary to assign a label **true** (t), **false** (f), or ϵ to each edge. In particular, a normal node has one outgoing edge labeled ϵ , whereas a predicate node has two outgoing edges corresponding to a true or false evaluation of the predicate. As an example, the CFG for bubble sort is shown in Figure 1(a). In our context, CFG is a standard code representation in reverse engineering to aid in understanding the structure of binary. However, while CFGs expose the control flow of a given code, they fail to provide other useful information, such as the way registers are manipulated by the code and the interaction between different functions.

2.2. Register Flow Graph

A Register Flow Graph (RFG) is used to capture how registers are manipulated by binary code, which is originally designed for authorship identification of binary code [1]. RFGs describe the flow and dependencies between registers as an important semantic aspect of the

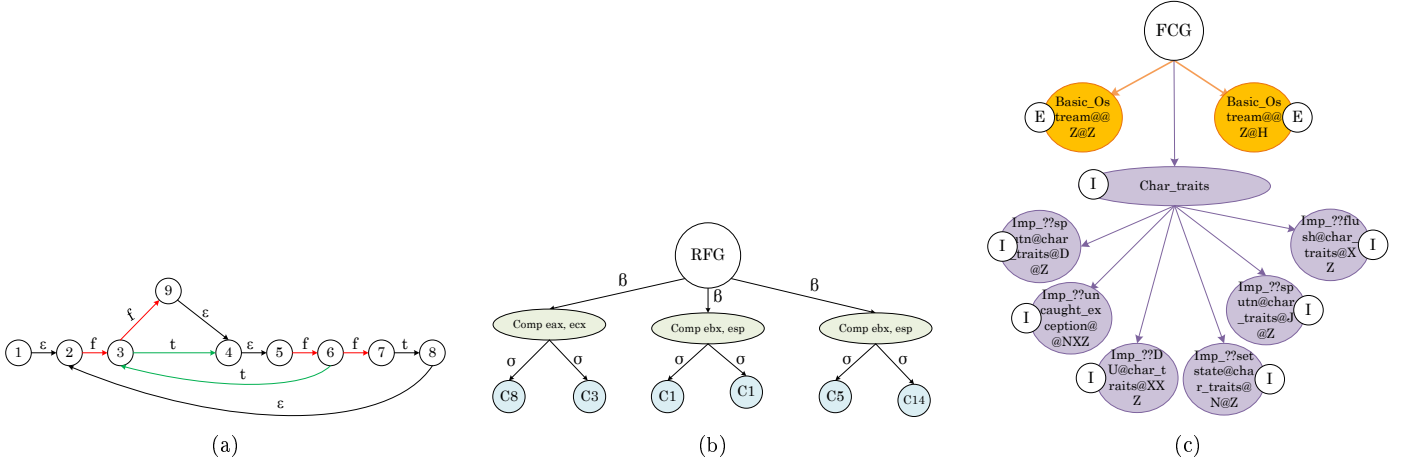


Figure 1: Classical representations for bubble sort function: (a) Control Flow Graph (b) Register Flow Graph (c) Function Call Graph

behavior of a program, which might indicate authorship as well as functionality. We briefly review the concept through an example shown in Figure 1(b). In the RFG, two labels are assigned to edges; β represents the basic block to which the compare instruction belongs (basic block id), and σ is the cost that is assigned based on the flow of the register values (instruction counts). Regardless of the number or complexity degree of functions, the following registers are often accessed: `ebp`, `esp`, `esi`, `edx`, `eax`, and `ecx`. Therefore, the steps involved in constructing a RFG for these registers are as follows.

- Counting the number of compare instructions,
- Checking the registers for each compare instruction,
- Checking the flow of each register from the beginning until the compare is reached,
- Classifying the register changes according to the 15 proposed classes in [1].

In RFGs, assembly instructions are classified into four families: Stack, Arithmetic, Logical, and Generic operation, as detailed in the following.

- Arithmetic: `add`, `sub`, `mul`, `div`, `imul`, `idiv`, etc.
- Logical: `or`, `and`, `xor`, `test`, `shl`.
- Generic: `mov`, `lea`, `call`, `jmp`, `jle`, etc.
- Stack: `push` and `pop`.

2.3. Function Call Graph

A Function Call Graph (FCG) is the representation of a function in binary code as a directed graph with labeled vertices, where the vertices correspond to functions and the edges to function calls. Two labels, I and E are assigned to the nodes; I represents internal library functions and E represents external library functions. An example of FCG for the bubble function is shown in Figure 1(c). In the literature, external call graphs have been used for malware detection [7]. In such a case, model graphs and data graphs are compared in order to distinguish call graphs representing benign programs from those based on malware samples [14, 7].

3. SIGMA Approach

In this section, we first provide an overview of the proposed SIGMA approach in Section 3.1. We then describe the three building blocks of an *SIG* in Section 3.2. We introduce the *SIG* concept in Section 3.3. Finally, we describe methods for *SIG* graph matching in Section 3.4.

3.1. Overview

The overall architecture of our SIGMA approach is depicted in Figure 2. There are two main phases: (i) training phase, and (ii) testing phase, detailed as follows.

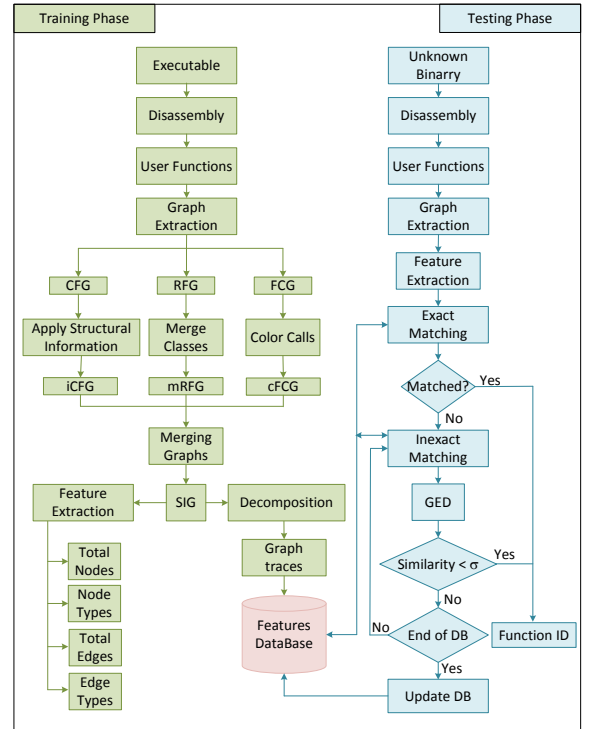


Figure 2: SIGMA Architecture

The training phase consists of four steps; (i) disassembling the executable and manually filtering out compiler-

related functions; (ii) constructing CFG, RFG, and FCG graphs from user functions; (iii) applying structural information to CFG to obtain the informational control flow graph, *iCFG*; applying new merged classes to RFG to obtain a merged register flow graph, *mRFG*; and applying colored classes to FCG to obtain a colored function call graph, *cFCG* (these concepts will be explained in Section 3.2). (iv) Merging the previous graphs into a single representation called *SIG*. We then decompose the *SIG* into a set of traces aiming to apply inexact matching between different graphs. Moreover, we consider various properties of the *SIG*, such as the total number of nodes, node types (data, control, dependence, or structural), edge types, total number of edges, the depth of the graph, etc. We save these details into a database with the function ID. On the other hand, given a set of unknown assembly instructions, the testing phase construct the *SIG* and extract the properties of the constructed graph and compare it with the existing *SIGs* graphs in the database. Hence, we have two methods for matching graphs: (i) *exact matching*: two graphs are said to match exactly if they have the same properties. (ii) *inexact matching*: it is based on edit distance calculation and the result is compared to predefined threshold value δ . Two functions are the same if their similarity score is less than δ . More formally, we have following definitions.

Definition 1. Let f_1, f_2 be two functions, we say f_1 is the copy (or origin) of f_2 , if $SIG(f_1)$ matches $SIG(f_2)$.

Definition 2. Let f_1, f_2 be two functions, and $SIG(f_1) \rightarrow a$ and $SIG(f_2) \rightarrow b$ denote extracting *SIG* traces a and b from f_1 and f_2 . Let $sim(a, b)$ be a similarity function and δ a predefined threshold value ($\delta < 1$). We say f_1 and f_2 are similar if $sim(a, b) < \delta$.

3.2. Building Blocks

In this section, we extend the existing representations introduced in Section 2 to form the building blocks of *SIG*.

3.2.1. Structural Information Control Flow Graph (*iCFG*)

As mentioned in Section 2, traditional CFGs consist of basic blocks each of which is a sequence of instructions terminating with a branch instruction. We can thus only obtain the structure of a function from a CFG. The lack of more detailed information in CFGs means two entirely different functions may yield the same CFG, which will cause confusion for identifying similar functions. Therefore, we extend standard CFGs with a colored scheme based on structural information about the probable role or functionality of each node. For example, if the majority of instructions in one node is arithmetic or logical, it may provide hints about the functionality of the node (e.g., cryptographic function usually involves a large number of for loops). By enriching standard CFGs with such information as different colors of nodes, which we call *iCFG*,

we have a better chance to distinguish two functions even if they have the same CFG structure. Table 1 shows some example categories of structural information we consider in coloring the nodes.

Table 1: Structural information categories

Category	Description
Data Transfer (DT)	Data transfer instructions such as mov, movzx, movsx
Test(T)	Test instructions such as cmp, test
ArLo	Arithmetic and logical instructions such as add, sub, mul, div, imul, idiv, and, or, xor, sar, shr
CaLe	System call, API call, and Load effective instructions such as lea
Stack	Stack instructions such as push, pop

The assignment of classes depends on two percentages: (i) the two highest percentages, and (ii) the lowest percentage, among the proposed categories. By considering the highest percentages, we aim to measure the majority category in the function. We choose two highest percentages because we have noticed that some classes, such as **Data Transfer**, are always dominant in many cases such that considering in addition the second highest percentage would provide more reliable coloring. Table 2 shows color classes for *iCFG*. Each row shows three classes. For example, the second row shows classes 1, 2, and 3; class 2 occurs when the two majorities are DT, T and the minority is Stack.

Table 2: Color classes for *iCFG*

Color Classes	Majority	Minority
1/2/3	DT, T	ArLo/Stack/CaLe
4/5/6	DT, ArLo	T/CaLe/Stack
7/8/9	DT, CaLe	ArLo/Stack/T
10/11/12	DT, Stack	T/CaLe/ArLo
13/14/15	T, ArLo	DT/CaLe/Stack
16/17/18	T, CaLe	DT/ArLo/Stack
19/20/21	T, Stack	DT/ArLo/CaLe
22/23/24	ArLo, Stack	T/DT/CaLe
25/26/27	ArLo, CaLe	Stack/DT/T
28/29/30	Stack, CaLe	T/DT/ArLo

As an example, by applying the color classes in Table 2 to Figure 1(a), we can obtain the *iCFG* shown in Figure 3(a). This *iCFG* involves five color classes: 22, 4, 3, 10, and 1. From Table 2, we can see that the majority of those classes belong to: **ArLo-Stack**, **DT-ArLo**, **DT-T**, **DT-Stack**, **DT-T**. This is reasonable since the main functionality of the bubble sort algorithm is manipulating values in an array and consequently the main action is transferring the values from one location to another, which explains the large number of DT instructions. As demonstrated by the example, by using this extended control flow graph *iCFG*, we can capture more semantic information that might be helpful in identifying functions in

binary code. Nonetheless, the *iCFG* only contains control information about basic blocks, and it lacks other useful semantics, such as the way registers are manipulated and the way functions interact with each other. Hence, we introduce two other building blocks in addition to *iCFG*.

3.2.2. Merged Register Flow Graph (*mRFG*)

As mentioned in Section 2, *RFG* is a binary code representation for capturing program behaviors based on an important semantics of the code, i.e., how registers are manipulated. The original *RFG* is designed for authorship attribution purposes, therefore it lacks support for some cases that are important for function identification: i) when both operands of `cmp` are constants (C), ii) when one of the operands is a constant and the other one is a register (`reg`), iii) when both operands are memory locations (ML), iv) when one of the operands is a memory location and the other one is a register, and v) when the operands are a mixture of constants and memory locations. These cases are especially important for identifying functions in binary code, and hence we extend the *RFG* by adding several new classes as shown in Table 3.

Table 3: Updated classes of register access

Class	Arithmetic	Logical	Generic	Stack	C C	C Reg	ML ML	ML Reg	ML C
1	1	0	0	0	0	0	0	0	0
2	0	1	0	0	0	0	0	0	0
3	0	0	1	0	0	0	0	0	0
4	0	0	0	1	0	0	0	0	0
5	1	1	0	0	0	0	0	0	0
6	1	0	1	0	0	0	0	0	0
7	1	0	0	1	0	0	0	0	0
8	1	1	1	0	0	0	0	0	0
9	1	1	0	1	0	0	0	0	0
10	1	0	1	0	0	0	0	0	0
11	1	1	1	1	0	0	0	0	0
12	0	1	1	0	0	0	0	0	0
13	0	1	0	1	0	0	0	0	0
14	0	0	1	1	0	0	0	0	0
15	0	1	1	1	0	0	0	0	0
16	0	0	0	0	1	0	0	0	0
17	0	0	0	0	0	1	0	0	0
18	0	0	0	0	0	0	1	0	0
19	0	0	0	0	0	0	0	1	0
20	0	0	0	0	0	0	0	0	1

Moreover, as another improvement over the original *RFG* representation, we merge certain nodes inside an *RFG*, e.g., class one and class two together are equivalent to class five. In this manner, we can reduce the number of nodes to improve the efficiency in analyzing an *RFG*. Finally, since the original *RFG* depends only on the `cmp` instructions, we also extend *RFG* instructions to the `test` instruction. After applying those extensions and modifications, we obtain the new representation *mRFG*, as shown in Figure 3(b). The *mRFG* has three more nodes than

its corresponding *RFG*; one of these is `test` instruction, and the other two are related to the immediate memory address and the constant. Moreover, we have merged the original classes (nodes in blue): *C8–C3* to *C8*, *C1–C1* to *C1*, and *C5–C14* to *C11*. The reference to the new classes *C17* and *C19* may provide useful semantics about the functions, for instance bubble sort function mainly deals with constants and sorts them in memory locations.

3.2.3. Color Function Call Graph (*cFCG*)

As mentioned in Section 2, traditional FCGs represent system calls in a binary code. Among a set of system calls $C = \{C_1, C_2, \dots, C_n\}$, each call may be either internal or external. To distinguish these from each other, we extend FCGs with a color scheme as follows. The function of coloring nodes defines the color class α in two cases. For an internal call, we only need one color, because internal system calls are mostly related to compiler functions rather than to user functions. As to external calls, we define the color classes using a range of values $0 < \alpha < 1$, because we may have various external system calls potentially connecting to API that is very important for identifying functions. More precisely, we extend FCGs to a new representation which we call *cFCG*, using the color function defined as follows.

$$f(c) = \begin{cases} \alpha = 0, & \text{if } c \text{ is internal system call} \\ 0 < \alpha < 1, & \text{if } c \text{ is external system call} \end{cases}$$

As an example, having applied this new representation to our running example, we obtain the *cFCG* shown in Figure 3(c). Besides serving as a building block of our proposed approach, the *cFCG* representation may also be helpful in other related tasks by highlighting the difference between various types of calls. This kind of graph may use for malware analysis, for instance, malware clustering through clustering external system calls.

3.3. SIG: Semantic Integrated Graph

The building blocks introduced in the previous section provide complementary views on binary code by emphasizing on different aspects of the underlying function semantics. Inspired by a recent work [19], in which different representations of source code are combined for vulnerability detection in source code (which is a different problem from ours as binary code lack much of the useful information available in source code), we combine those different but complementary representations of binary code into a joint data structure in order to facilitate more efficient graph matching between different binary code for identifying reused functions. Formally, a semantic integrated graph (*SIG*) is defined as follows.

Definition 3. A semantic integrated graph is a directed attributed graph $G = (N, V, \zeta, \gamma, \vartheta, \lambda, \gamma, \psi, \omega, \delta)$ where N is a set of nodes, $V \subseteq (N \times N)$ is a set of edges and γ is

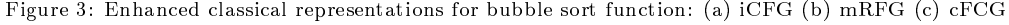


Table 4: Some traces of bubble sort function *SIG*

Node	Traces	Traces Type
22	ϵ , ϵ , C8 t	out in
4	f, C14 ϵ , f	out in
22 OR 4	ϵ , C8, C14, f f, t, ϵ	out in
4 AND 3	f f	out in

Table 5: Graph Features for bubble sort function

Features	Frequency
Total # of Nodes	15
Total # of Edges	18
# of Control Nodes	5
# of Control Edges	8
# of Call Nodes	4
# of Register Nodes	6
Connected Graphs	3
k-Cone	1,2

for this purpose. The edit distance between two graphs measures their similarity in terms of the number of edits needed to transform one into the other [10]. We implement this concept as follows. Given two *SIGs*, we define the following two elementary traces to transform one graph into another: **Edge-edit traces**, including ϱ_{κ_r} , re-labels the edge, and **Node-edit traces**, including ϱ_{ν_r} , re-colors the node by merging nodes from the other graph into one node. An edit edge $V_{G,H}$ between two *SIGs* G and H , is defined as a set of sequence of traces $(\varrho_1, \varrho_2, \dots, \varrho_n)$ such that $G = \varrho_n (\varrho_1 (\varrho_{n-1}(H) (\dots \varrho_1(H) \dots)))$. To quantify this similarity, the weight of all edit traces is measured, i.e., $V = (\varrho_n, \varrho_2, \dots, \varrho_n)$ as $w(V) = \sum_{i=1}^n w(\varrho_i)$. The edit distance between two *SIGs* is thus defined as the minimum weight of all edit edges and nodes between them, i.e., $sim(G, H) = \min w(V_{G,H})$. The distance measure between the nodes follows the same reasoning, with operations instead of traces. In Algorithm 1, we calculate the graph edit distance between two *SIGs* G and H , by measuring the cost of transforming G to H . The algorithm starts by extracting the traces of the two graphs as mentioned earlier, and then checks the cost of transforming each node in G to nodes in H , and finally calculates the total cost.

We define the dissimilarity between two *SIGs* G and H as follows.

Definition 5. The dissimilarity $\rho(G, H)$ between two *SIGs* is a value in $[0, 1]$, where 0 indicates the graphs are the most similar and 1 indicates the least similar, as formulated in the following.

$$\rho(G, H) = \frac{w(V_{G,H})}{|N_G| + |N_H| + |V_G| + |V_H| + |\varrho_G| + |\varrho_H|}$$

Algorithm 1: Graph Edit Distance

Input: G, H : Semantic Integrated Graphs (*SIG*)

Output: $sim(G, H)$: Similarity result between two *SIGs*

Initialization;

$\varrho_g, \varrho_h \leftarrow \emptyset;$ // traces

$w(V_{G,H}) \leftarrow 0;$ // weighted cost

$sim(G, H) \leftarrow 0;$ // similarity score

begin

$\varrho_g \leftarrow extract_traces(G);$

$\varrho_h \leftarrow extract_traces(H);$

foreach ϱ_i in ϱ_g **do**

foreach ϱ_j in ϱ_h **do**

$w_j(V_{G,H}) = cost_of_transforming(\varrho_i, \varrho_j);$

$sim_{i,j}(G, H) = \min w(V_{G,H});$

$sim(G, H) = sim(G, H) + sim_{i,j}(G, H);$

return $sim(G, H)/i;$

Where $w(V_{G,H})$ is the weighted cost of traces, $|N_G|$ and $|N_H|$ are the number of nodes, $|V_G|$ and $|V_H|$ are the number of edges, and $|\varrho_G|$ and $|\varrho_H|$ are the number of traces in G and H respectively.

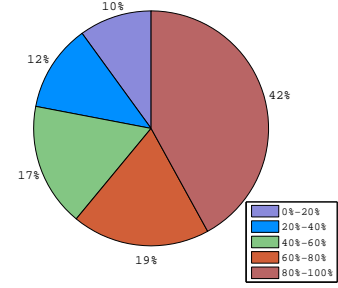


Figure 6: Similarity scores of function variants

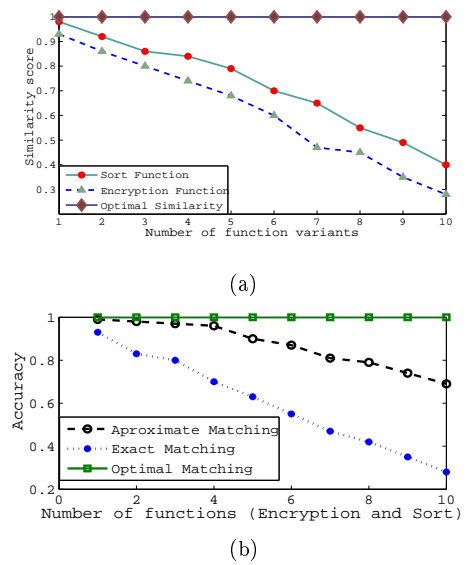


Figure 7: (a) The relation between the number of variants and the similarity score (b) The accuracy of using exact and approximate matching

shown in Figure 9, which captures the similarity in terms of traces. Based on the *SIGs* and their traces, we conduct exact matching and inexact matching, whose results are shown in Table 9 and Table 10, respectively. The exact matching shows that the similarity is 78.5%, and the cost for inexact matching is 0.311 (close to 0).

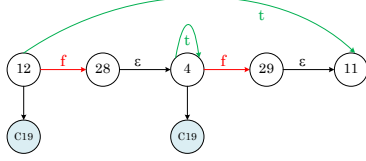


Figure 9: *SIG* for RC4 in Zeus

Table 9: Exact Matching between *SIG*-RC4 in Citadel and Zeus

Features	<i>SIG</i> -RC4 (Zeus)	<i>SIG</i> -RC4 (Citadel)	Similarity
Total # of Nodes	7	9	78%
Total # of Edges	8	11	72%
# of Control Nodes	5	5	100%
# of Control Edges	6	7	86%
# of Call Nodes	1	0	50%
# of Register Nodes	2	3	67%
Connected Graphs	3	3	100%
K-Cone	1,2,3,4	1,2,3	75%
Average Similarity			78.5%

Table 10: Inexact Matching between *SGF*-RC4 in Citadel and Zeus

Citadel Node	Zeus Node	Costs	Node(s) with Minimum Cost	Cost %
13	12	1 out, 0 in	12	(1/10)
	28	1 out, 1 in	-	
	4	1 out, 2 in	-	
	29	1 out, 1 in	-	
	11	0 out, 2 in	-	
28	12	3 out, 0 in	-	0
	28	0 out, 0 in	28 (Select this)	
	4	3 out, 1 in	-	
	29	0 out, 0 in	29	
	11	0 out, 1 in	-	
4	12	0 out, 0 in	12 (Already chosen)	0
	28	1 out, 1 in	-	
	4	0 out, 0 in	4 (Select this)	
	29	1 out, 1 in	-	
	11	0 out, 0 in	11	
11	12	2 out, 0 in	-	(1/10)
	28	1 out, 0 in	28 (Already chosen)	
	4	2 out, 2 in	-	
	29	1 out, 0 in	29 (Select this)	
	11	0 out, 2 in	-	
12	12	2 out, 0 in	-	(1/9)
	28	1 out, 1 in	-	
	4	2 out, 0 in	-	
	29	1 out, 1 in	-	
	11	0 out, 1 in	11	
Total Cost				0.311

6. Limitations and Future Direction

The previous section shows that our proposed *SIGMA* approach yields promising results for identifying reused functions in binary code. Nonetheless, the approach still has following limitations which we would like to address in our future work.

- Like most existing approaches, *SIGMA* assumes that binary code is already de-obfuscated. We note that, however, *SIGMA* can in fact address certain forms of obfuscation, such as register reassignments and code recording. Our future work will investigate this potential direction.
- As a learning-based approach, *SIGMA* also requires training data of known functionalities with multiple variants in order to collect sufficient features prior to its application to given code. To this end, we intend to build a feature database for common functionalities by applying *SIGMA* to publicly available code.
- We have not investigated the impact of different compilers in this study. However, we believe that *SIGMA* can overcome some of the changes caused by compilers with the rich set of structural information it captures, and we will confirm this in future work.
- Our future work will also evaluate the capability of the proposed method for dealing with fragments of functions.

7. Related Work

Our main inspiration comes from the recent work of combining different source code representations for discovering vulnerabilities [19]. In addition to the idea of combining different sources of information, we also borrow the definitions of graph traces. However, we note that the authors deal with a very different problem, which is to model and detect vulnerabilities, than ours, which is to identify reused functions. Another major difference is that we work on binary code instead of source code. This implies that we must employ entirely different representations, since much of the useful information available in source code, such as abstract syntax trees, is not applicable to binary code. To the best of our knowledge, this is the first effort on the use of multiple sources of structural information for binary code analysis. For identifying reused functions in binary code, existing work mainly fall into two categories: (i) static, and (ii) dynamic. We briefly review static approaches since our work is based on static analysis. Within static analysis, there exist some work that employ graphical representations. Interprocedural control flow using the call graphs of a program have been compared to show similarity to existing malware [10], where common nodes between two program call graphs are discovered. The authors in [12] build their graph by transforming a portable executable into a call graph with nodes and edges that capture system calls and system call sequences, respectively. They then convert the graph to a code graph to expedite analysis. The authors in [7]

propose a method where each malware sample is represented as an API call graph by integrating API calls and operating system resources to represent graph nodes. The authors in [17] compare metric values and introduce transformers and formulas that could use training data to generate a measure of the similarities between two procedures in binary code. The authors in [18] propose a method to identify malware variants based on a function-call graph. The authors in [8] develop a pragmatic effective code size reduction technique that exploits the structural similarity of functions. Unlike most existing work, our approach employs multiple sources of structural information to define the distance between variants of functions, which allows for more reliable and efficient matching. Our approach also differs from many existing work in the capabilities of inexact matching and matching function fragments based on graph traces.

8. Conclusion

The reverse engineering of binary code is an important but challenging task that demands automated techniques for preprocessing and cleaning the code. The identification of reused functions in binary code is one of the important aspect of this issue that has received limited attention in comparison with other aspects of binary analysis. In this paper, we have presented a novel approach called *SIGMA* for effectively identifying reused functions in binary code. Instead of relying on one source of information, our approach combines multiple representations into one joint data structure *SIG*. *SIGMA* also supports inexact matching and exact matching based on traces of the *SIG* which deals with function fragments. Both experimental results and case study have demonstrated the effectiveness of our method, and we have described several potential improvements to the approach in the previous section.

References

- [1] Alrabae, S., Saleem, N., Preda, S., Wang, L., Debbabi, M., OBA2: An Onion approach to Binary code Authorship Attribution. Digital Investigation, 2014, Vol 11, PP. S94-S103, Elsevier.
- [2] Balliu, M., Dam, M., Guanciale, R., Automating Information Flow Analysis of Low Level Code, In Proceedings of the 21st ACM Conference on Computer and Communication Security, (CCS'14), 2014, ACM.
- [3] Bencsáth, B., Buttyán, L., Félégyházi, M., Pék, G. sKyWIper (aka Flame aka Flamer): A complex malware for targeted attacks, 2012.
- [4] Calvet, J., Fernandez, J. M., Marion, J. Y., Aligot: cryptographic function identification in obfuscated binary programs, In Proceedings of the 2012 ACM conference on Computer and communications security, 2012, PP. 169-182, ACM.
- [5] Cesare, S., Xiang, Y., Zhou, W., MalwiseŪan effective and efficient classification system for packed and polymorphic malware, IEEE Transcation on Computers, 2013, Vol 62, PP. 1193-1206, IEEE.
- [6] David, Y., Yahav, E., Tracelet-based code search in executables, In Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation, 2014, PP. 349-360, ACM.
- [7] Elhadi, A. A. E., Maarof, M. A., Barry, B. I., Hamza, H., Enhancing the Detection of Metamorphic Malware using Call Graphs, Computers and Security Journal, 2014, Vol 46, PP. 62-78, Elsevier.
- [8] Edler von Koch, T. J., Franke, B., Bhandarkar, P., Dasgupta, A. Exploiting function similarity for code size reduction, In Proceedings of the 2014 SIGPLAN/SIGBED conference on Languages, compilers and tools for embedded systems, 2014, PP. 85-94, ACM.
- [9] Gröbert, F., Willems, C., Holz, T., Automated identification of cryptographic primitives in binary programs, In Recent Advances in Intrusion Detection, 2011, PP. 41-60, Springer Berlin Heidelberg.
- [10] Hu, X., Chiueh, T. C., Shin, K. G., Large-scale malware indexing using function-call graphs. In Proceedings of the 16th ACM conference on Computer and communications security, 2009, PP. 611-620, ACM.
- [11] Khoo, W. M., Mycroft, A., Anderson, R., Rendezvous: A search engine for binary code, In Proceedings of the 10th Working Conference on Mining Software Repositories, 2013, PP. 329-338, IEEE Press.
- [12] Lee, J., Jeong, K., Lee, H., Detecting metamorphic malwares using code graphs, In Proceedings of the 2010 ACM symposium on applied computing, 2010, PP. 1970-1977, ACM.
- [13] Myles, G., Collberg, C., K-gram based software birthmarks, In Proceedings of the 2005 ACM symposium on Applied computing, SAC '05, PP. 314-318, ACM.
- [14] Riesen, K., Jiang, X., Bunke, H., Exact and inexact graph-matching: methodology and applications, Managing and Mining Graph Data Advances in Database Systems, 2010, Vol 40, PP. 217-247, Springer.
- [15] Ruttenberg, B., Miles, C., Kellogg, L., Notani, V., Howard, M., LeDoux, C., Pfeffer, A., Identifying Shared Software Components to Support Malware Forensics, In Detection of Intrusions and Malware, and Vulnerability Assessment, 2014, PP. 21-40, Springer International Publishing.
- [16] Sæbjørnsen, A., Willcock, J., Panas, T., Quinlan, D., Su, Z., Detecting code clones in binary executables, In Proceedings of the eighteenth international symposium on Software testing and analysis, 2009, PP. 117-128, ACM.
- [17] Stojanovic, S., Radivojevic, Z., Cvetanovic, M., Approach for Estimating Similarity between Procedures in Differently Compiled Binaries, Information and Software Technology, 2014, Elsevier.
- [18] Xu, M., Wu, L., Qi, S., Xu, J., Zhang, H., Ren, Y., Zheng, N., A similarity metric method of obfuscated malware using function-call graph, Journal of Computer Virology and Hacking Techniques, 2013, Vol 9, No 1, PP. 35-47, Springer.
- [19] Yamaguchi, F., Golde, N., Arp, D., Rieck, K., Modeling and Discovering Vulnerabilities with Code Property Graphs, Proc. of 35th IEEE Symposium on Security and Privacy, 2014, IEEE.