

Universitatea POLITEHNICA din București

Facultatea de Automatică și Calculatoare,
Catedra de Calculatoare



LUCRARE DE DIPLOMĂ

Analiza aplicațiilor de tip malware

Conducător Științific:
As.dr.ing. Laura Gheorghe

Autor:
Cristian Condurache

București, 2013

University POLITEHNICA of Bucharest

Automatic Control and Computers Faculty,
Computer Science and Engineering Department



BACHELOR THESIS

Malware Analysis

Scientific Adviser:

As.dr.ing. Laura Gheorghe

Author:

Cristian Condurache

Bucharest, 2013

TODO:

Write acknowledgements

Maecenas elementum venenatis dui, sit amet
vehicula ipsum molestie vitae. Sed porttitor
urna vel ipsum tincidunt venenatis. Aenean
adipiscing porttitor nibh a ultricies. Curabitur
vehicula semper lacus a rutrum.

Quisque ac feugiat libero. Fusce dui tortor,
luctus a convallis sed, lacinia sed ligula.
Integer arcu metus, lacinia vitae posuere ut,
tempor ut ante.

Abstract

Malware is currently a major security threat for computers and smartphones, with efforts being taken into improving malware detectors with behavior-based detection. In order to classify applications, malware detectors need some form of malicious behavior specification which are usually identified manually by researchers. We present a Linux implementation of the malspec-mining algorithm which automates this process. This algorithm recognizes such specifications by comparing known malicious and benign applications. The output consists of behavior patterns which are specific to the inputted malware and that do not occur in benign applications.

Keywords: behavior-based detection; malspec-mining algorithm; malicious behavior; kernel programming

Contents

Acknowledgements	i
Abstract	ii
1 Introduction	1
2 State of the Art	3
2.1 Malware types	3
2.2 Avoiding detection	4
2.3 Malware detectors	5
2.4 Malspec-Mining Algorithm	5
3 Malspec Mining Algorithm and Design	7
3.1 Mining Minimal Contrast Subgraph Patterns	7
3.2 Maximal Common Edge Set	8
3.3 Malspec Mining Algorithm	8
3.4 Design	9
3.4.1 System Call Interceptor Driver	9
3.4.2 Network Interceptor	9
3.4.3 Obtaining traces	10
3.4.4 Computing specifications	11
4 Implementation	12
4.1 System Call Interceptor Driver	12
4.2 Network Interceptor	13
4.3 Reading traces	14
4.4 Building dependence graphs	15
4.5 Malspec algorithm	15
5 Evaluation	17
5.1 Malware test environment	17
5.2 Test scenario	17
5.3 Statistics	19
6 Conclusions and Future Work	22
A Configuring Monitored System Calls	23
A.1 syscalls.xml	23
B Malsharp output sample	25
B.1 program_test, diff_test	25

List of Figures

2.1	Malware Statistics	4
3.1	System call dependence graph	9
3.2	Architecture	10
5.1	Virtual Machine test configuration	17
5.2	Maximal common edge set. (a) - program_test; (b) - diff_test.	18
5.3	Minimal transversal of complement on the malware graph. (a) - program_test; (b) - diff_test.	19

List of Tables

5.1	System calls for program_test and diff_test	18
5.2	Malware sample analysis for android.xml	20
5.3	Malware sample analysis for all_syscalls.xml	20

Notations and Abbreviations

API – Application Programming Interface
DDoS – Distributed Denial of Service
DOM – Document Object Model
fd – file descriptor
ICMP – Internet Control Message Protocol
malspec – malware specification
NAT – Network Address Translation
NI – Network Interceptor
pid – process id
procfs – proc filesystem
SB – signature database
SCID – System Call Interceptor Driver
symlink – symbolic link
syscall – system call
TCP – Transmission Control Protocol
UDP – User Datagram Protocol
VM – Virtual Machine
XML – Extensible Markup Language

Chapter 1

Introduction

From large corporations to the average user, computer and network environment security is an important requirement to which malware is a threat. Malicious software is a program that has been written by an attacker to fulfill a harmful intent. In order to achieve this, the program has to interact with the victim's operating system..

The number of users of a specific operating system is directly correlated to the degree of interest malware writers take in developing software to target that specific operating system. Due to Linux's increasing popularity, better security for operating systems that are based on the Linux kernel has become a necessity. This supports the need for developing tools for Linux malware analysis and improving malware detection methods.

Earlier detection methods focused on analyzing the contents of the executable file of the malware program, such as identifying instruction sequences which were characteristic for specific malware instances. These methods performed poorly when confronted with unknown malware or new variants of existing ones. Also, in response, attackers started to write malware that modifies its own file while replicating itself, thus eluding these detection methods.

This resulted in a switch to developing behavior based detection systems that are independent from the exact contents of the executable file. Therefore, when analyzing malware samples, analysts started to search for program behavior patterns that suggest a malicious intent. In order for these patterns to work, programs need a higher-level common behavior specification.

The system call interface meets this requirement as malware needs to interact with the operating system to achieve its goals and it is common to all malware. A typical malware example would be an executable file that replicates itself by reading its own file and then copying it to system directories. This can be captured in a behavior pattern which, compiled into malware specifications, can then be used by malware detectors in order to classify programs based on their behavior.

The project presented in this thesis, named Malsharp, is a Linux tool for automatically searching for malicious program behavior patterns. This tool is a Linux implementation of the malware specification mining algorithm which identifies these behavior patterns by comparing known malware samples to known benign programs. These patterns are a collection of Linux system call parameter dependencies that capture the malicious behavior.

Malsharp is intended to be used by malware specialists to help them analyze new malware samples. It can also be used as part of an automatic detection mechanism which classifies programs based on their behavior by using the malicious behavior specifications.

This thesis organized as follows: in [Chapter 2](#) we describe the main types of malware threats, current detection types, methods that malware uses to avoid detection and a short presentation

of the malware specification mining algorithm.

In [Chapter 3](#) we define several notions used in the algorithm, a detailed view of the malspec-mining algorithm and we reveal our own architecture. [Chapter 4](#) presents the implementation details for each element of architecture from kernel modules to userspace application.

[Chapter 5](#) describes our test environment, explains a short test scenario and how the algorithm works and shows the results we obtained with our implementation on actual malware samples. [Chapter 6](#) presents the conclusions and future work.

Chapter 2

State of the Art

Computer security mostly refers to the mechanisms that are used to protect computers and networks from different threats. Recent history is full of famous attacks on computers, networks and especially over the Internet. Although this field deals with many more security related issues, one of the major threats to computer security is malware.

2.1 Malware types

Malware, or malicious software, is software programmed and used by attackers in order to gain access to private computers, to obtain sensitive information or to simply disrupt normal computer operation. Malware generically refers to a variety of program forms: viruses, worms, Trojan horses and spyware.

A virus is a program that attempts to replicate itself into other executable files by injecting or replacing code. When the infected programs are run, they can infect other ones in turn. They typically target common programs which are found on most machines and executables and copy themselves to key directories.

Worms are similar to viruses, only that they spread over the network to other hosts which have the same vulnerability as the initial host. This is done by performing network port scans, DNS queries and then trying to infect other machines. Also, this type of malware usually downloads a second program from a remote server.

Trojan horses, or Trojans, are non-self-replicating code that try to gain privileged access to an operating system and while they seem to be performing a legitimate action they deliver a malicious payload. The payload often contains a backdoor for the attacker that gives him access to the computer or a botnet to send spam or perform Denial-of-service attacks.

A botnet is a network of “zombie” computers that have been compromised and are used to perform malicious actions like: distributed denial-of-service attacks (or DDoS). DDoS attacks are carried out in order to stop web services from functioning efficiently or from functioning at all.

Adware, or advertising-supported software, is a term that refers both legitimate advertising software and malware. When used to refer malware, adware presents unwanted advertisements to the user of a computer, sometimes under the form of pop-up windows.

Spyware is software created for gathering information without the user’s consent. This kind of malware is usually installed without the knowledge of the user or by using deceptive tactics [12].

In Figure 2.1 we can see the percentage of malware samples that correspond to the each type. These statistics were produced by **Panda Security** on April 16th, 2011.

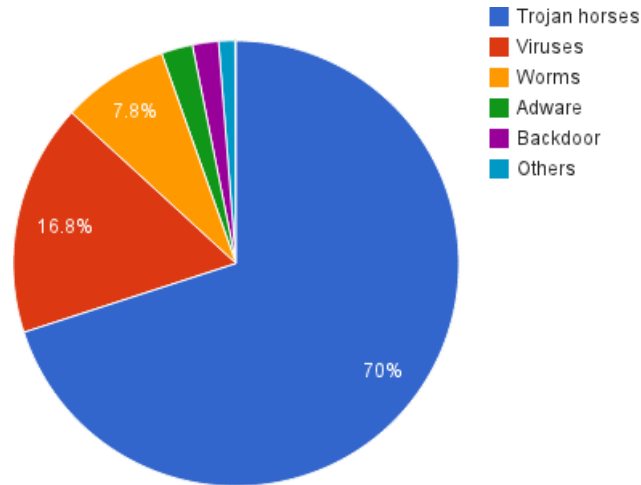


Figure 2.1: Malware Statistics

2.2 Avoiding detection

Malware and antimalware software evolution is tightly coupled: when detection methods become more efficient, malware writers use better hiding techniques. In response, anti-malware developers use better algorithms. In the course of time, the development of antimalware software has led attackers to start using different obfuscation methods to avoid detection, such as: concealing API-calling behavior, polymorphic and metamorphic malware.

Malware writers obscure API-calling behavior by using indirect calls. In this respect, API calls can be achieved by using hard coded addresses but this method is incompatible with different versions of the operating system. Another method commonly employed is to define homonymy functions that first locate the API functions addresses and then use the stored address. Also, arrays can be used to store the API function name and the relocation address [8].

Polymorphic viruses use a polymorphic engine to change its executable while keeping the original function intact. A common technique to “morph” viruses is to encrypt the malicious payload and decrypt it at runtime. The encrypted code will then appear to be meaningless and it will be ignored. To obfuscate the decryption routine, the code is transformed by inserting **nop** instructions, permuting register allocation, reordering instructions and inserting jump instructions to maintain the semantics.

Metamorphic viruses also use an engine to change their code in a variety of ways, such as code transposition and substituting instruction sequences with equivalent ones. In addition to this, they interweave their code with the original program’s code to trick heuristic detection methods. An important difference between a metamorphic engine and a polymorphic one is that the first can rewrite itself while the second cannot [6].

2.3 Malware detectors

Malware detectors are developed by performing analysis on samples gathered through various means: honeypots, web crawlers, spam traps and security analysts that collect them from infected computers. Bayer *et al.* [4] provided insight into common behavior by analyzing almost one million malware samples by monitoring their network activity and tracking data flows.

Bayer *et al.* created a platform, named Anubis [2], for the dynamic analysis of malware samples which targeted Windows operating systems. The behavior of malware samples was monitored for file system, registry, network and botnet activity, GUI windows and sandbox detection. Sandboxes are contained environments used to run and test malicious software. The statistics they presented offer insight into common malware behavior and give a hint to what the main goal of malware detectors should be.

Signature based malware detectors use a list of signatures (signature database or SB) to identify known viruses. The signatures are computed by applying a hash function on the malware file. If a part of a program matches a signature entry from the list, then it is classified as malware. This detection method performs very poorly when confronted with new samples because the signature is unknown. Also, malware writers can easily avoid detection from this type of detectors by using obfuscation techniques in their programs, like polymorphism or metamorphism [7].

Over time, the approach in detecting malware has evolved from analyzing the contents of infected executable files towards identifying malicious or potentially malicious behavior patterns. These patterns are extracted from the malware sample by static or dynamic analysis.

Static analysis of the executable involves scanning the file for particular instruction sequences or different API calls. In order to avoid detection from this type of analysis, attackers attempt to obscure their API-calling behavior or they use a polymorphic engine.

Another method for analysis is to monitor the behavior of the malicious program during runtime in a sandbox, otherwise known as dynamic analysis. This method monitors the malware's interaction with the operating system and the network traffic it produces in order to determine its behavior.

Semantics-aware malware detectors can overcome the problems posed by obfuscation by using specifications of malicious behavior which are not affected by polymorphic malware. By using a higher-level specification, different versions or implementations of malware which perform the same behavior can be detected. Another advantage of this type of detector is that it can also successfully classify unknown malware [10].

2.4 Malspec-Mining Algorithm

The problem with behavior-based detection is that the required specifications have to be manually identified by a malware specialist. The -mining algorithm developed by Christodorescu *et al.* [7] provides a method for automating this otherwise time consuming task.

Their malspec-mining algorithm starts by collecting execution traces from malware and benign programs, then it constructs the corresponding dependence graphs and then it computes the specification of malicious behavior as difference of dependence graphs as minimal contrast subgraph patterns [11].

The malspec-mining algorithm was implemented and tested on a Windows operating system and, although it identified a large number of malware specifications, it also managed to capture most of the specifications that were indicated by specialists [7].

In this paper we present an implementation of this algorithm for GNU/Linux based operating systems. In order to capture a program's behavior we developed a system call interceptor and a network traffic interceptor as Linux kernel modules.

Then, a user space program reads the traces from the kernel module and constructs a graph where each node represents a system call and the edges represent parameter dependencies. The edges are determined by interpreting the parameter type, direction and value of the recorded system calls. Finally, the malspec-mining algorithm will generate the malicious behavior specifications.

Chapter 3

Malspec Mining Algorithm and Design

3.1 Mining Minimal Contrast Subgraph Patterns

The minimal contrast subgraph patterns used in the malspec mining algorithm were introduced by R. Ming Hieng Ting *et al.* [11]. The following definitions will provide a better understanding of the malspec mining algorithm.

An **edge set** is a labelled graph, *i.e.* a graph with labels attached to its vertices and edges, that has no isolated vertices.

Given two graphs, G_p and G_n , C is a **common edge set** if and only if C is an edge set and it is a common subgraph. C is a **maximal common common edge set** if it is a common edge set of the two graphs and if and only if there does not exist a superset which is also a common edge set.

The notions of **maximal common edge set** and **minimal contrast edge sets** are connected and we can determine the second set by obtaining the complement of the first with respect to the original graph and then computing the minimal transversal.

Given G_p and $\{G_{n1}, G_{n2}, \dots, G_{nk}\}$, let M_i be the set of maximal common edge sets between G_p and G_{ni} . Then the set of minimal transversals of $\bar{M}_1 \cup \bar{M}_2 \cup \dots \cup \bar{M}_k$ will be the set of all **minimal contrast edge sets** between G_p and $\{G_{n1}, G_{n2}, \dots, G_{nk}\}$, where \bar{M}_i is the graph complement of M_i with respect to G_p .

For a graph, a **partition** is a set of disjoint and not empty subsets, named cells, of V , *i.e.* the set of all the vertices in the graph. All the vertices in the same cell have the same label and vertices from different cells are labelled differently. The union of all the cells in the partition is equal to V .

Given G_p and G_n that are associated with the partitions T_p and T_n , a **minimal contrast vertex set** is a subset of a cell from T_p such that its cardinality is larger by 1 in comparison with the cells from T_p .

The **minimal contrast subgraph** of a positive graph, G_p , with respect to a negative graph, G_{ni} , is the minimal union of the **minimal contrast edge sets** and the **minimal contrast vertex sets** of the two graphs. The minimal union works like normal union but it removes any graphs that are supergraphs of others in the set.

3.2 Maximal Common Edge Set

In determining the minimal contrast edge set, the most demanding computational task is finding the maximal common edge set. The problem of testing whether a subgraph relationship exists between two graphs is NP-complete.

The maximal common edge set is determined using the backtrack algorithm proposed by J. McGregor for computing the maximal common subgraph [9].

The algorithm takes as input two graphs, G_1 and G_2 , with $|V_1| < |V_2|$, and for each node in the G_1 it tries to find a correspondent node from G_2 while maximizing the number of matching edges. The result of the algorithm is a list of pairs of corresponding nodes.

In order to maximize the number of matching edges of the two graphs, a mapping matrix MARCS is used to indicate if an arc from G_1 can correspond to G_2 . Initially, all the edges from the first graph can correspond to any edge from the second graph. When a node is paired with another, the edges connected to it can only correspond to edges that are connected to the other node. The edges that do not correspond are marked in MARCS with a zero value.

3.3 Malspec Mining Algorithm

The malspec-mining algorithm collects the execution traces and uses them to create a dependence graph for each program.

Each node of the dependence graph represents one system call with its arguments, while edges represent dependences between arguments of different system calls. The dependence graphs are constructed after aggregating similar operations like multiple reads and writes using the same buffer and file handle, thus resulting in fewer nodes.

Edges are established between nodes which present def-use dependencies: if a previous (in execution order) system call parameter has an out (or inout) parameter with the same value as an in (or inout) argument of a later system call then an edge can be established from the first node to the second node.

Then, the minimal contrast subgraph miner operates in three stages. First, the maximal common edge set for the malware graph and all the benign graphs is determined using the algorithm designed by McGregor [9].

Next, the common edge sets are unioned together and the minimal traversals of their complements are computed, thus obtaining the minimal contrast edge sets for the malware graph.

Finally, the contrasts are minimally unioned with the minimal contrast vertex sets to give the complete set of minimal contrast subgraphs.

The resulting subgraphs for each comparison of the malware sample with a benign program are maximally unioned, *i.e.* removing graphs which are subgraphs of others, giving the desired malicious behavior specifications.

The following example shows how an [execution trace](#) is transformed into a [dependence graph](#). Every system call from [Listing 3.1](#) is represented by a node in the dependency graph as seen in [Figure 3.1](#).

Then, each node is compared only with the nodes that succeed it in execution order. In this example, the three system calls are connected two def-use dependency edges because they operate on the same file descriptor.

```

1 open("/bin/ls", O_RDWR) = 3;
2 read(3, 0x80000001, 255) = 127;
```



```
3 close(3) = 0;
```

Listing 3.1: System call trace

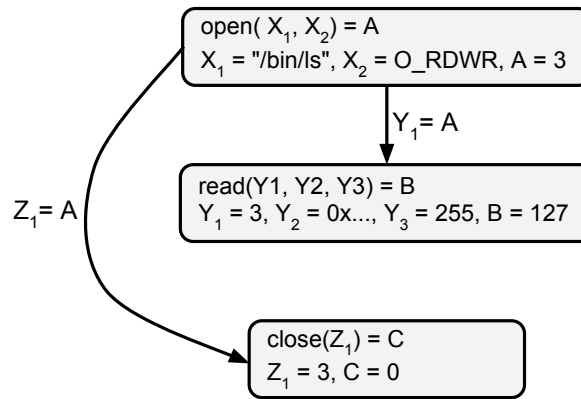


Figure 3.1: System call dependence graph

3.4 Design

Malsharp takes as input three pathfiles to:

- an XML file which contains the system calls, with argument type and direction, that will be considered for the analysis,
- the malware sample that will be analyzed and
- a file which contains a set of benign programs, one per line.

The architecture of Malsharp that was used in implementing the algorithm is presented in [Figure 3.2](#).

3.4.1 System Call Interceptor Driver

The system call interceptor is a kernel module used for obtaining execution traces for specific processes. It can be controlled from userspace through different `ioctl` commands such as: setting the process id to be monitored and setting the system calls to be monitored. Also, each execution trace entry for a system call can be read and removed, the trace history can be cleared and the total number of entries can be read.

This module will be configured to monitor each malicious or benign program from userspace by our application. The resulting execution trace will then be read and used to generate the dependency graph.

3.4.2 Network Interceptor

The network interceptor is a kernel module that was created for monitoring network traffic for specific protocols, such as: TCP, UDP and ICMP. It is also controlled through `ioctl` calls and

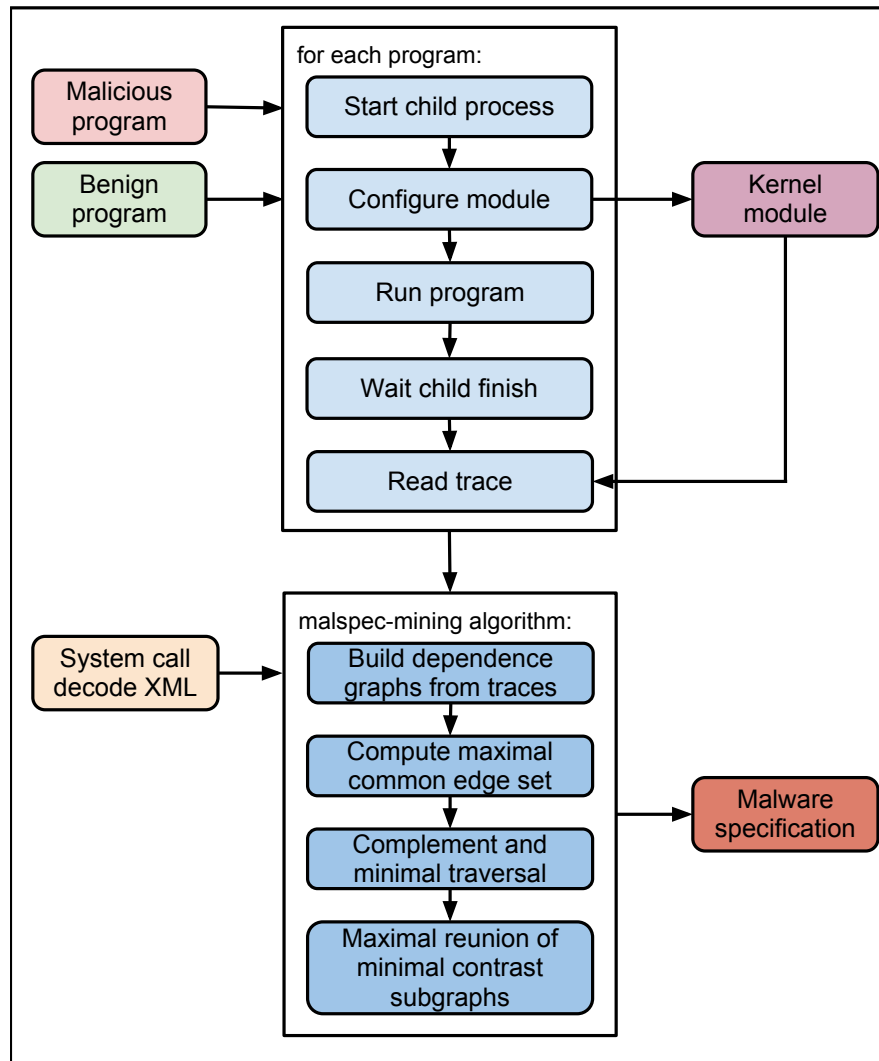


Figure 3.2: Architecture

it can be configured to monitor traffic by inspecting source or destination ports for TCP and UDP.

3.4.3 Obtaining traces

The malspec-mining algorithm is run in user space and it will create a new process for each program that it will analyze. Before running the program it will configure the kernel module to monitor the new process and then they will record its execution trace. Only system calls that have entries in the XML input file will be monitored.

The main program will wait for the child process to finish and then it will start reading the execution trace from the system call interceptor driver.

3.4.4 Computing specifications

The malspec-mining algorithm receives the execution trace as input and uses it to create the dependence graph for each program. Because the information gained from monitoring the system calls contains only the register values, the type and direction information for each argument must be filled. This is done by parsing the XML file given as input to Malsharp.

Next, each pair of graphs is run through the different stages of the algorithm. The output consists of a set of malware specifications that are behavior patterns found in the malware sample that did not occur in any of the benign programs. These specifications can later be used for detecting malware.

Chapter 4

Implementation

Malsharp was implemented in C++, a natural choice considering we had to communicate with the Linux kernel modules which are implemented in C. Both the kernel modules and the userspace application use the same header file that contains the common data structures required for configuring the modules and reading the execution traces.

4.1 System Call Interceptor Driver

The System Call Interceptor Driver (SCID) can be configured by giving the pid of the process as an argument at module insertion or by opening the device and calling `ioctl`. The system calls to be monitored can only be configured by `ioctl`.

The SCID registers itself as a character device named **scid** using the misc device interface provided by the Linux kernel. It also creates a symlink to `procfs` that is used by the userspace part of the application for opening the device.

This driver logs only the system call number and the arguments and return values. The system calls are intercepted by replacing the original value of entry in the system call table with the address of an interceptor function which will run the system call and it will perform the necessary logging.

System call history is implemented with a queue, which contains the recorded system call values and the actual system call number. When an `ioctl` read call is received from userspace, it will remove and return the first entry in the queue. It is also possible to clear the entire log history if needed or to get the total number of entries from the queue.

The kernel module uses the `syscall_params` struct to retrieve the parameters from the stack when intercepting system calls. The data structure returned by the kernel is called `sctrace_t` and it contains the system call number, a `syscall_params` struct with the argument values and the return value. The following listing contains the definitions for the data structures used by the kernel module.

```
1  /* data struct for syscall intercepting */
2  struct syscall_params {
3      long ebx, ecx, edx, esi, edi, ebp, eax;
4  };
5
6  /* data struct for read syscall history */
7  typedef struct _sctrace_t {
```

```

8      int sc_no;
9      struct syscall_params sc_params;
10     long ret;
11 } sctrace_t;

```

Listing 4.1: SCID data structures

The SCID module `ioctl` commands are defined using the macros provided by the Linux kernel. For the necessary code value required in the kernel macros we chose `0xA1`, which is currently unused according to the documentation file in the kernel.

```

1 #define IOCTL_SET_PID          _IOW(0xA1, 1, int)
2 #define IOCTL_ADD_SYSCALL      _IOW(0xA1, 2, int)
3 #define IOCTL_DEL_SYSCALL      _IOW(0xA1, 3, int)
4 #define IOCTL_CLEAR_HISTORY    _IO(0xA1, 4)
5 #define IOCTL_COUNT_HISTORY    _IOR(0xA1, 5, int)
6 #define IOCTL_READ_HISTORY     _IOR(0xA1, 6, sctrace_t)

```

Listing 4.2: Command macros for `ioctl`

- `IOCTL_SET_PID` sets the pid of the process that will be monitored
- `IOCTL_ADD_SYSCALL` start monitoring a system call
- `IOCTL_DEL_SYSCALL` stop monitoring a system call
- `IOCTL_CLEAR_HISTORY` deletes all entries with trace information
- `IOCTL_COUNT_HISTORY` returns the number of trace entries
- `IOCTL_READ_HISTORY` get and remove the first trace entry

4.2 Network Interceptor

The Network Interceptor (NI) is also a kernel module that can be configured in a similar manner to the SCID. This module uses kernel netfilter hooks to intercept packets on `PRE_ROUTING` chain and it logs statistics such as total number of packets and total size transmitted.

The NI is controlled through `ioctl` calls which receive a `intercept_info_t` struct containing the protocol, source and destination ports it should monitor. A zero value destination or source port enables monitoring for all the traffic.

```

1 typedef struct _intercept_info_t {
2     /* IPPROTO_*, TCP, UDP, ICMP */
3     unsigned char xport_protocol;
4     /* used only in TCP and UDP */
5     unsigned short int source;
6     unsigned short int dest;
7 } intercept_info_t;

```

Listing 4.3: parameter data structures

Like SCID, the network interceptor also makes a symlink to `procs` from which statistics can be read. The information that is collected includes the total number of packages transmitted and the total information size sent for each pair of source and destination ports.

The information gathered by the NI is not used in the determining the malware specifications at the moment, but integration is possible. Currently, the module can be used to verify if a sample malware program tried to use the network to spread itself or to send information.

4.3 Reading traces

The entries recorded by the SCID module do not contain any type and direction information about the system call. Therefore, in order to successfully identify def-use dependences, the type and direction have to be properly set in a `syscall_t` struct, which is defined below, before the graph is constructed.

```

1  /* param data structure to describe a syscall's parameters */
2  #define MAX_NUM_PARAM    8
3
4  typedef struct param {
5      unsigned char type; /* fd, int, unsigned int, char*, void*,
        unsigned short, size_t */
6      unsigned char dir; /* 1 in, 2 out, 3 inout */
7      long value;
8  } param_t;
9
10 typedef struct _syscall_t {
11     int syscall_no; /* the system call number */
12     param_t param[MAX_NUM_PARAM]; /* last element is the return
        value */
13 } syscall_t;

```

Listing 4.4: parameter data structures

The `syscall_t` struct contains a `param` array with a maximum of 8 argument entries. The first 7 entries are used for the register values of the system call and the 8th entry holds the return value. Possible argument directions are in, out and inout, while data types considered are int, unsigned int, file descriptor, char*, void*, unsigned short and size_t.

Malsharp keeps system call decoding information in an XML file which contains for each system call the type and direction information for each register. The file is parsed using the open source library named **pugixml** [3].

We chose this library because the XML configuration files do not contain many entries and **pugixml** uses DOM representation. Also, it supports XPath queries, which we used to query information about particular system calls when we added type and direction information to the execution trace.

Each program is run separately by creating a new process, making it yield the processor by waiting on a named semaphore and then replacing its image with `execve`. The new process has to be forced to yield the processor in order for the parent process to have sufficient time to set the child's process pid as the target for the interceptor modules. Otherwise, a partial system call trace would be obtained instead and the results would be compromised.

After the parent process has configured the kernel module to monitor its child process, it increments the semaphore to enable it to continue execution. After that, the parent process waits for the child process to end and then it starts reading the execution trace.

4.4 Building dependence graphs

The execution trace inputted is used to create a graph for the sample program. Each node of the graph will contain a `syscall_t` struct. On creation, each graph object receives an array of these structs from which to create the graph's nodes.

After the nodes are created, each pair of nodes is verified for def-use dependences and if such dependencies exist then a new edge is added in a vector as a pair of node pointers. To determine if a def-use dependency exists between two nodes, the two `param` arrays from the `syscall_t` structs are compared.

In addition to system call information, each node contains a unique identification number (or index) in order to uniquely identify nodes within the same graph, even if they have the same system call number and argument values.

Although file descriptors are returned as integer values in Linux, they are treated as different data types because integer values are a common argument type for system calls.

Treating file descriptors the same way as integers would result in false def-use dependences. If an open system call would return file descriptor number 3 and a read call on a different file would try to read 3 characters, then a false def-use dependence edge would be generated between this pair of nodes because the fd and the integer have the same value.

Also, due to the fact that Linux reuses file descriptors, an additional check had to be implemented when generating the dependencies. Otherwise, false extra def-use edges would be generated that do not reflect actual argument dependencies.

As pointed out by Christodorescu *et al.* in [7], consecutive system calls of the same type which operate on the same resources can be aggregated into a single node. On a file descriptor, one read of N bytes is equivalent to N consecutive reads of 1 byte.

Aggregation is particularly important because it reduces the number of nodes that have to be paired in the backtracking algorithm designed by McGregor for determining the maximal common edge set. According to [7], two nodes, n_1 and n_2 , should be aggregated if:

- n_1 and n_2 had the same system call
- n_1 and n_2 shared a common def-use predecessor, p
- there was no node between n_1 and n_2 that performed a different system call on the same resource.

In addition to these, we added an extra constraint:

- for each in-edge, or out-edge, that connects n_1 to another node n_3 , a corresponding in-edge, or an out-edge, respectively, that connects n_2 to n_3 must exist.

Although the fourth constraint we added reduces the number of aggregated nodes, it ensures that information related to consecutive nodes that have other def-use dependencies to other resources is not lost.

4.5 Malspec algorithm

McGregor's backtrack algorithm was implemented using recursion, although the original proposed pseudocode was iterative.

The initial algorithm assumed that the graphs are not labelled, but in [11] labelling is presented as one of several powerful pruning methods. Apart from aggregation, label pruning was also

used in implementing the backtracking algorithm: only the nodes which shared the same label, the system call number, were tentatively paired.

After the complements of the maximal common edge sets are computed, the minimal transversal of their union must be computed. In our implementation all the edges represent def-use dependencies, so there is no “weight” attached to the edge labels. Therefore, the minimal transversal of the graph can be computed by doing a simple breadth first search and keeping a list of the edges that added a new node to the queue.

Testing for subgraph and supergraph relationships between two graphs was performed by finding the maximal common subgraph and then testing if every node of a graph was matched by another node from the second graph.

Chapter 5

Evaluation

5.1 Malware test environment

Malsharp was tested in a sandbox environment: a VMware virtual machine with a distribution of Ubuntu with the Linux 3.7.8 kernel which has been bridged to the local network for Internet access.

To ensure that all the malware samples were tested in the same environment, we used snapshots so that each file was tested on the same virtual machine state. In this way, different malware samples could not influence each other's execution.

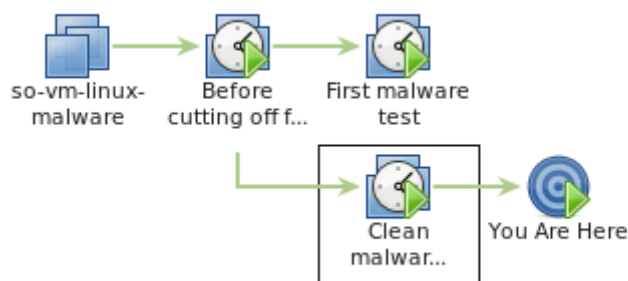


Figure 5.1: Virtual Machine test configuration

In [Figure 5.1](#) the snapshot test configuration is shown. The “Clean malware state” is the malware test snapshot we mentioned earlier. “Before cutting off for malware test” is the VM state before changing the connection type from NAT to bridge and before removing all ssh keys from the machine.

5.2 Test scenario

In a typical test scenario, Malsharp receives an XML file which contains the system calls with parameter information that we want to monitor in our executables. We can test for malware specifications only on a subset of the total system calls a Linux operating system has.

For a better comprehension on how the malspec mining algorithm works, we will consider the following example. Let us assume that we want to search for malware specifications with `open`, `read`, `write` and `close` system calls. In this example, `program_test` will be the malware sample and `diff_test` will be the benign program.

The following table presents the system calls that each executable file makes and the relevant argument information for the dependency edges.

Table 5.1: System calls for `program_test` and `diff_test`

	<code>program_test</code>	<code>diff_test</code>
1	<code>open(...)</code> = fd1	<code>open(...)</code> = fd1
2	<code>close(fd1)</code>	<code>close(fd1)</code>
3	<code>open(...)</code> = fd2	<code>open(...)</code> = fd2
4	<code>read(fd2, ...)</code>	<code>read(fd2, ...)</code>
5	<code>close(fd2)</code>	<code>close(fd2)</code>
6	<code>open(...)</code> = fd3	<code>open(...)</code> = fd3
7	<code>write(fd3, ...)</code>	<code>close(fd3)</code>
8	<code>write(fd3, ...)</code>	<code>open(...)</code> = fd4
9	<code>write(fd3, ...)</code>	<code>write(fd4, ...)</code>
10	<code>close(fd3)</code>	<code>read(fd4, ...)</code>
11	-	<code>close(fd4)</code>

In Figure 5.2 the generated graphs for the two programs presented in Table 5.1 are shown. For each node we have written its index in the graph and its label.

Although there are three `write` calls in `program_test` that operate on the same file descriptor, these nodes are not aggregated because of the fourth restriction we added to node aggregation in Section 4.4. The three nodes cannot be aggregated because of the edges that exist between them.

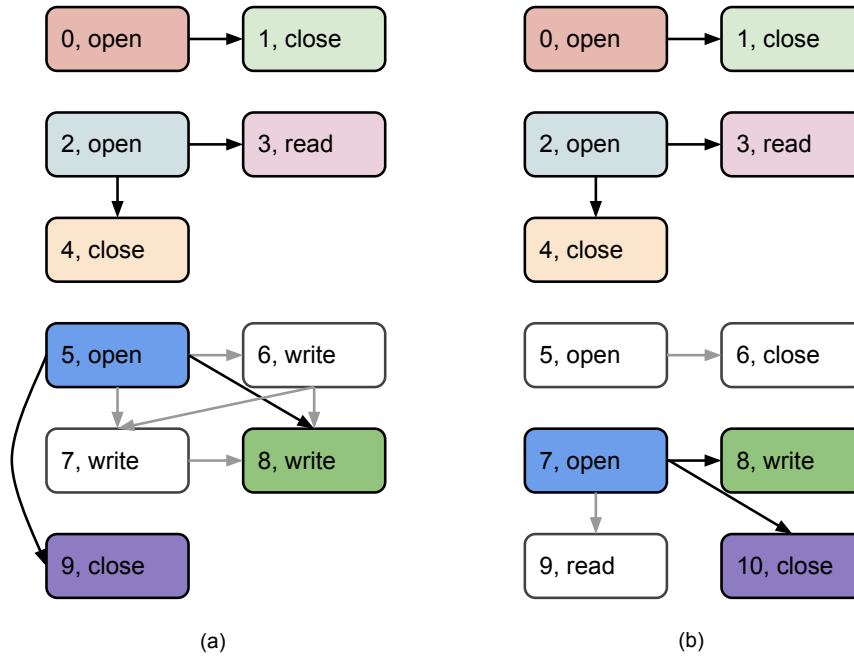


Figure 5.2: Maximal common edge set. (a) - `program_test`; (b) - `diff_test`.

The maximal common edge set that was determined using McGregor’s algorithm is shown by drawing paired nodes with the same colors. The nodes that do not have a matching node in the other graph have no color. The edges connecting them are gray because they are not a part of the maximal common edge set.

After computing the maximal common edge set, its complement in the malware graph is computed. In this particular example, the complement is a single connected graph, so the minimal union step of the algorithm is not required.

Then, a breadth first search is used to find the minimal transversal of the complement. The result of these two steps of the algorithm can be seen in [Figure 5.3](#). The malspec contains two nodes which are colored red and the black edge connecting them. Because we used a single beign program to find the malspec, this will be the final result of the algorithm.

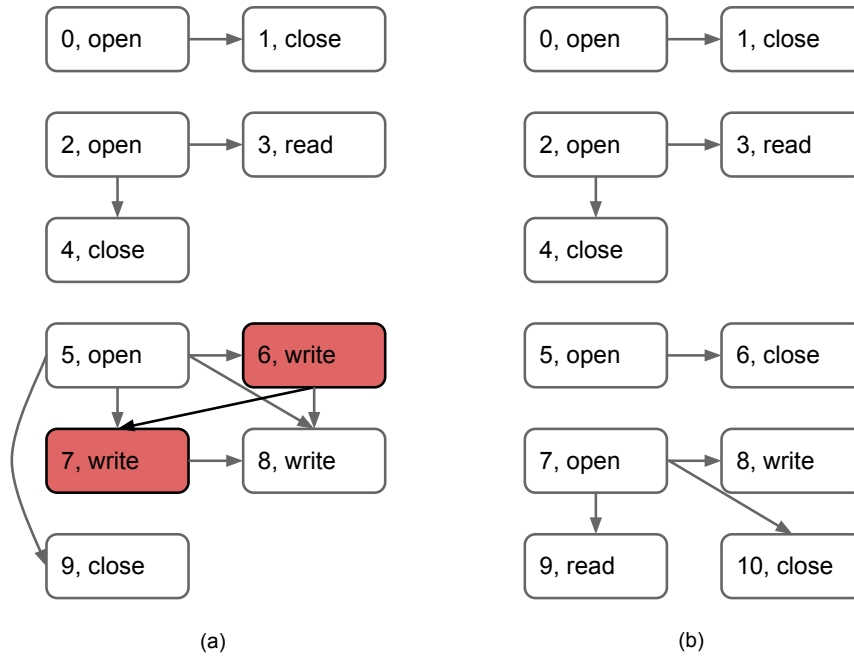


Figure 5.3: Minimal transversal of complement on the malware graph. (a) - `program_test`; (b) - `diff_test`.

The complete output from Malsharp for `program_text` and `diff_test` can be seen in [Appendix B.1](#).

5.3 Statistics

Malsharp was tested using the Linux malware samples provided by [1]. All the malware samples that were tested also had additional classification information provided by antivirus software.

As pointed out by Burguers *et al.* in [5], the increasing number of smartphones on the market with the Android platform, which has a Linux kernel, make malware analysis an important issue.

Therefore, we used two sets of system calls to monitor the malicious and benign programs: a smaller set which was indicated by Burguers *et al.* in their paper about behavior-based

detection on Android [5] and a larger set which contained other system calls that might be used for malicious purposes.

The larger set is a superset of the other and has more system calls that use with file descriptors and with socket communication. Only a part of this XML file has been added to this paper as [Appendix A.1](#) because the original file was too large.

The results for each set of system calls can be seen in the following two tables, one per set. The set of benign programs we considered for testing consisted of: `ls`, `lsmod`, `ping` and `cat`.

Table 5.2: Malware sample analysis for `android.xml`

Sample	N	N'	E	Malspecs	Time	Observed behavior
Backdoor.Linux.CGI.a	12	12	5	1	10.208s	open, read, close
Backdoor.Linux.Phobi.1	36	35	16	1	0.0641s	open, read, close
Trojan.Linux.Rootkit.n	12	11	5	1	0.558s	open, read, close
Virus.Linux.Osf.8759	692	442	285	0	0.725s	blocked driver
Virus.Linux.Radix	14	13	4	1	0.786s	open, read
Virus.Linux.Silvio.b	46	46	24	1	0.983s	open, read, close
Virus.Linux.Snoopy.c	128	93	39	1	1.271s	open, read, close
Virus.Linux.Svat.b	23	23	12	0	0.617s	-

In [Table 5.2](#), we show the results for 8 malware samples. Although the initial collection of malware samples had 20 samples, some samples could not be tested because the operating system on our VM did not have all the shared libraries it required in order to run. Also, other samples were actually tools for flooding a range of ip addresses in a network.

Most of the samples generated dependency graphs with a very small number of nodes. This is a consequence of the fact that `android.xml` contains a very small number of system calls used for monitoring: `read`, `open`, `close`, `chmod`, `lchown` and `access`.

`Virus.Linux.Osf.8759` generated the biggest graph and we can see that the node aggregation method reduced the total number of nodes considerably, from 692 to 442. Unfortunately, this malware sample interfered with our SCID module, and analysis was not possible.

Table 5.3: Malware sample analysis for `all_syscalls.xml`

Sample	N	N'	E	Malspecs	Time	Observed behavior
Backdoor.Linux.CGI.a	13	13	5	1	0.827s	open, read, close
Backdoor.Linux.Phobi.1	36	35	16	1	0.641s	open, read, close
Trojan.Linux.Rootkit.n	14	14	5	1	1.044s	open, read, close
Virus.Linux.Osf.8759	19	19	9	0	53m 31s	open, read, fstat read, close; fork
Virus.Linux.Radix	25	22	6	1	0.673s	open, read, write; creat, write
Virus.Linux.Svat.b	25	25	12	0	3.495s	replaces <code>stdio.h</code>

Analysis results for the larger set of system calls is shown in [Table 5.3](#). The total analysis took longer due to the fact that the generated dependence graphs are larger, especially in the case of the benign programs.

The most remarkable result was for the `Virus.Linux.Osf.8759` sample, although it had the longest analysis time - 53 minutes and 31 seconds. The malspec it returned consisted of two graphs.

The first graph contains the following nodes: `open`, `read`, `fstat`, `read` and `close`, all connected by file descriptor dependencies. The second graph contains a single node, a `fork`

system call. The child process was not monitored.

Another interesting malware sample was `Virus.Linux.Radix` which inserted itself into the executable file for the `ls` command, `/bin/ls`, and `cp` command, `/bin/cp`. It also searched the current directory for executable files and it infected those as well.

Chapter 6

Conclusions and Future Work

Computer security is an important field and continuous advances ensure protection from threats such as disruption of normal computer and network operation, information leaks and attackers who benefit financially from compromised hosts.

As Linux based operating systems are becoming evermore popular, malware security for these systems will continue to grow in importance. Also, use in embedded systems makes Linux an important target for malware writers.

The **Malsharp** tool represents a proof of concept for a Linux implementation of the malware specification mining algorithm and leaves room for developing a more advanced detector. There are still a lot of possible improvements that are planned for implementation.

First, the traces are currently obtained by using the system call interceptor driver. In the future we plan to switch to using strace instead. Also, using iptables instead of the network interceptor module might be a possible improvement.

Secondly, another possible improvement includes adding other dependency edges apart from def-use, such as: comparing strings for common substrings or self-referential dependencies for when the malware file opens its own file.

Finally, the backtracking algorithm for determining the maximal common subgraph can be further improved with other pruning methods like node ordering strategies and ancestor-leaf equivalence pruning.

Malsharp will be released under an Open Source license to permit others to contribute to this project. We hope that with time, this will become a useful and reliable tool for Linux malware analysts.

Appendix A

Configuring Monitored System Calls

A.1 syscalls.xml

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <syscalls>
3     <syscall no="2" name="fork">
4         <ebx type="int" dir="in"/>
5     </syscall>
6     <syscall no="3" name="read">
7         <ebx type="fd" dir="in"/>
8         <ecx type="pvoid" dir="out"/>
9     </syscall>
10    <syscall no="4" name="write">
11        <ebx type="fd" dir="in"/>
12        <ecx type="pvoid" dir="in"/>
13    </syscall>
14    <syscall no="5" name="open">
15        <ebx type="pchar" dir="in"/>
16        <ecx type="int" dir="in"/>
17        <edx type="ushort" dir="in"/>
18        <ret type="fd" dir="out"/>
19    </syscall>
20    <syscall no="6" name="close">
21        <ebx type="fd" dir="in"/>
22        <ret type="void" dir="none"/>
23    </syscall>
24    <syscall no="8" name="creat">
25        <ebx type="pchar" dir="in"/>
26        <ecx type="ushort" dir="in"/>
27        <ret type="fd" dir="out"/>
28    </syscall>
29    <syscall no="9" name="link">
30        <ebx type="pchar" dir="in"/>
31        <ecx type="pchar" dir="in"/>
32    </syscall>
```

```

33     <syscall no="10" name="unlink">
34         <ebx type="pchar" dir="in"/>
35     </syscall>
36     <syscall no="11" name="execve">
37         <ebx type="pchar" dir="in"/>
38         <ecx type="pchar" dir="in"/>
39         <edx type="pchar" dir="in"/>
40     </syscall>
41     <syscall no="12" name="chdir">
42         <ebx type="pchar" dir="in"/>
43     </syscall>
44     <syscall no="15" name="chmod">
45         <ebx type="pchar" dir="in"/>
46         <ecx type="ushort" dir="in"/>
47     </syscall>
48     <syscall no="16" name="lchown">
49         <ebx type="pchar" dir="in"/>
50         <ecx type="ushort" dir="in"/>
51         <edx type="ushort" dir="in"/>
52     </syscall>
53     <syscall no="19" name="lseek">
54         <ebx type="fd" dir="in"/>
55         <ecx type="int" dir="in"/>
56         <edx type="int" dir="in"/>
57     </syscall>
58     <syscall no="20" name="getpid">
59         <ret type="int" dir="out"/>
60     </syscall>
61     <syscall no="23" name="setuid">
62         <ebx type="uint" dir="in"/>
63     </syscall>
64     <syscall no="24" name="getuid">
65         <ret type="uint" dir="out"/>
66     </syscall>
67     <syscall no="33" name="access">
68         <ebx type="pchar" dir="in"/>
69         <ecx type="int" dir="in"/>
70     </syscall>
71     <syscall no="37" name="kill">
72         <ebx type="uint" dir="in"/>
73         <ecx type="int" dir="in"/>
74     </syscall>
75     <syscall no="38" name="rename">
76         <ebx type="pchar" dir="in"/>
77         <ecx type="pchar" dir="in"/>
78     </syscall>
79     <syscall no="39" name="mkdir">
80         <ebx type="pchar" dir="in"/>
81         <ebx type="ushort" dir="in"/>
82     </syscall>
83 </syscalls>

```

Listing A.1: System Calls Decode Information XML (syscalls.xml)

Appendix B

Malsharp output sample

B.1 program_test, diff_test

```
1 Monitored syscalls are:
2 3 4 5 6
3 Getting system call trace for ../test/program_test
4 driver contains 10 syscall entries
5 Malware graph:V[10]={ (0 sc=5), (1 sc=6), (2 sc=5), (3 sc=3), (4 sc
   =6), (5 sc=5), (6 sc=4), (7 sc=4), (8 sc=4), (9 sc=6), };
6 E[10]={ (0, 1), (2, 3), (2, 4), (5, 6), (5, 7), (5, 8), (5, 9), (6,
   7), (6, 8), (7, 8), };
7
8 Analyzing malware vs ../test/diff_test
9 Getting system call trace for ../test/diff_test
10 driver contains 11 syscall entries
11 Benign graph:V[11]={ (0 sc=5), (1 sc=6), (2 sc=5), (3 sc=3), (4 sc
   =6), (5 sc=5), (6 sc=6), (7 sc=5), (8 sc=4), (9 sc=3), (10 sc=6),
   };
12 E[7]={ (0, 1), (2, 3), (2, 4), (5, 6), (7, 8), (7, 9), (7, 10), };
13
14 -1 -1 -1 -1 -1 0 -1 -1 -1 1
15 -1 -1 -1 -1 -1 7 -1 -1 8 10
16 -1 -1 0 -1 1 7 -1 -1 8 10
17 -1 -1 2 3 4 7 -1 -1 8 10
18 0 1 2 3 4 7 -1 -1 8 10
19
20 Malspecs found:
21 Malspec 0
22 V[2]={ (0 sc=4), (1 sc=4), };
23 E[1]={ (0, 1), };
```

Listing B.1: Output for comparing program_test and diff_test

Bibliography

- [1] Open malware. <http://oc.gtisc.gatech.edu:8080/search.cgi?search=linux>.
- [2] Anubis. <http://anubis.isecclab.org>, 2007.
- [3] pugixml, light-weight c++ xml library. <http://pugixml.org>, 2012.
- [4] U. Bayer, I. Habibi, D. Balzarotti, E. Kirda, and C. Kruegel. A view on current malware behaviors. LEET'09 Proceedings of the 2nd USENIX conference on Large-scale exploits and emergent threats: botnets, spyware, worms and more, April 2009.
- [5] I. Burguera, U. Zurutuza, and S. Nadjm-Tehrani. Crowdroid: Behavior-based detection system for android. SPSM '11 Proceedings of the 1st ACM workshop on Security and privacy in smartphones and mobile devices, October 2011.
- [6] M. Christodorescu and S. Jha. Testing malware detectors. Proceedings of the ACM SIGSOFT International Symposium on Software Testing and Analysis, July 2004.
- [7] M. Christodorescu, S. Jha, and C. Kruegel. Mining specifications of malicious behavior. ESEC-FSE'07 Proceedings of the 6th joint meeting of the European software engineering conference and the ACM SIGSOFT symposium on The foundations of software engineering, September 2007.
- [8] W. Fu, J. Pang, R. Zhao, Y. Zhang, and B. Wei. Static detection of api-calling behavior from malicious binary executables. International Conference of Computer and Electrical Engineering, December 2008.
- [9] James J. McGregor. Backtrack search algorithms and the maximal common subgraph problem. Software - Practice and Experience, vol. 12, 23-34, 1982.
- [10] M.D. Preda, M. Christodorescu, S. Jha, and S. Debray. A semantics based approach to malware detection. 34th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, January 2007.
- [11] R. Ming Hieng Ting and J. Bailey. Mining minimal contrast subgraph patterns. 6th SIAM international conference on Data Mining, 2006.
- [12] M.F. Zolkipli and A. Jantan. Malware behavior analysis: Learning and understanding current malware threats. Second International Conference on Network Applications, Protocols and Services, September 2010.