# Malware Classification and Composition Analysis: A Survey of Recent Developments

Adel Abusitta, Miles Q. Li and Benjamin C. M. Fung

McGill University, Montreal, Canada

Email: adel.abusitta@mcgill.ca, miles.qi.li@mail.mcgill.ca,
ben.fung@mcgill.ca

**Abstract**

Malware detection and classification are becoming more and more challenging, given the complexity of malware design and the recent advancement of communication and computing infrastructure. The existing malware classification approaches enable reverse engineers to better understand their patterns and categorizations, and to cope with their evolution. Moreover, new compositions analysis methods have been proposed to analyze malware samples with the goal of gaining deeper insight on their functionalities and behaviours. This, in turn, helps reverse engineers discern the intent of a malware sample and understand the attackers' objectives. This survey classifies and compares the main findings in malware classification and composition analyses. We also discuss malware evasion techniques and feature extraction methods. Besides, we characterize each reviewed paper on the basis of both algorithms and features used, and highlight its strengths and limitations. We furthermore present issues, challenges, and future research directions related to malware analysis.

*Keywords:* Malware analysis, Malware classification, Security, anti-analysis techniques, Composition analysis

## 1. Introduction

In the recent years, many cyber-security mechanisms have been designed and developed to defend against evolving security threats. Nevertheless,

---

*Corresponding Author

recent statistics [1] indicate that malware are still evolving and becoming more sophisticated than ever. As a result, they become harder to detect and understand their innerworkings. This mainly stems from two essential reasons. The first is that attackers have now become more proficient in launching attacks and hiding their malicious behavior using anti-analysis techniques such as obfuscation and packing. The second reason is that the current communication and computing infrastructure is becoming more and more dynamic and heterogeneous, which enables a single malware to take various forms that are semantically but not structurally similar. This, in turn, makes malware analysis even more challenging.

Malware (or Malicious software) is a software that is designed to harm users, organizations, and telecommunication and computer system. More specifically, malware can block internet connection, corrupt an operating system, steal a user's password and other private information, and/or encrypt important documents on a computer and demand ransom. For the latest years, malware has been a growing threat to computer users and in 2017 the number of new malware increased by 22,9% over 2016 to reach 8,400,058 [2, 3, 4, 5]. Moreover, malware has become the primary medium to launch large-scale attacks, such as compromising computers, bringing down hosts and servers, sending out spam emails, crippling critical infrastructures and penetrating data centers [6, 7, 8]. These attacks lead to severe damage and significant financial loss [9, 10, 11].

Most antivirus engines detect and classify malware by continuously scanning files and comparing their signatures with known malware signatures. The malware signatures are typically created by human antivirus experts (known as malware defenders) who examine the collected malware samples. These malware signatures can be filename, text strings, or regular expressions of byte code [12, 13]. Obviously, signature-based methods can only detect traditional malware that do not change significantly. However, malware can hide its malicious behavior using anti-analysis techniques such as obfuscation, packing, polymorphism and metamorphism, in such a way that the code would look quite different from its original version. Thus, the primary shortcoming of the signature-based method is that they entail high precision but low recall. Also, the process of creating malware signatures is labor-intensive. Considering that there is a large number of new malware that appear every day, there is a pressing need to develop new intelligent malware analysis methods to tackle the challenges.

To alleviate the burden of manual signature crafting, researchers propose

automatic signature generation methods [14, 15]. The content of the signatures can be Windows system call combinations [16], control flow graph [15], and functions [14].

Researchers also propose to use machine learning models to detect and classify malware [12, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27]. Different from other machine learning-driven classification tasks, such as image classification, there is a competition between malware creators and defenders. When malware defenders propose a new malware analysis system using some features and machine learning models, malware creators often update their malware design to avoid being detected. Then malware defenders would propose new systems to detect and analyze the new generation of malware and so forth. The race between malware defenders and attackers may never come to an end.

Recently, many researchers have started to use deep learning models to enhance the detection and classification accuracy of malware classification [24, 25, 26, 27]. Although promising results have been achieved through the ability to extract robust and useful features using the state-of-the-art deep learning architectures, the proposed models were shown to be highly vulnerable to adversarial examples, which can be easily designed (simply by perpetuating parts of the inputs) by attackers to fool Artificial Intelligence (AI)-driven malware analysis systems and make them generate erroneous decisions [24, 25, 26, 27, 28, 29]. As a result, several methods have been proposed to defend against adversarial examples [28, 29].

In addition to malware classification, researchers in malware analysis have improved new techniques and methods to analyze the composition of malware samples by matching their functionalities and behaviours to multiple known malware families. This, in turn, helps reverse engineers discern the intent of a malware sample and the attacker. Moreover, these composition methods enable the reverse engineers and organizations to effectively triage their resources.

## 1.1. The Scope

This literature review classifies and compares the recent and main findings in malware classification. Unlike other similar works which only focus either on AI-driven malware classification [30] [31] [32] or on non-AI-driven malware classification [33] [34], this paper includes both AI-driven and non-AI-driven recent works. We are also surveying methods and approaches that

recently have been proposed to analyze the composition of malware samples, in order to understand their functionalities and behaviours. To the best of our knowledge, this is the first work that survey the existing composition analysis techniques. This survey also aims at identifying the main issues and challenges related to recent malware classification and composition analysis techniques. In particular, our analysis leads to recognize three major problems to address. The first is the need to overcome modern evading techniques (or anti-analysis techniques) such as metamorphism. The second relates to the efficiency and scalability of malware search engines as the number of functions in the repository might need to scale up to millions. The third concerns the vulnerability of malware classification system to evolving adversarial examples. We also uncover possible topics that need further study and investigation, such as sustainable malware analysis system. In this regard, we propose a few guidelines to prepare efficient and trustworthy malware detection and analysis system.

### 1.2. Contribution

The main contributions of this survey are:

- Proposing a new taxonomy for describing and comparing the recent and main findings in malware classification and composition analysis.

- Designing a new framework for analysing the existing malware classification and composition analysis techniques.

- Identifying and presenting open issues and challenges related to malware analysis.

- Identifying a number of trends on the topic, with guidelines on how to improve existing solutions to address new and continuing challenges.

### 1.3. Organization

The rest of this paper is organized as follows. In Section 2, we discuss the related survey papers. In Section 3 and Section 4, we present the proposed taxonomy for organizing reviewed malware classification and composition analysis approaches, respectively. Section 5 characterises reviewed papers according to the proposed taxonomy. The challenges and current issues are pointed out in Section 6. Section 7 suggests possible research topics in malware analysis. Finally, Section 8 concludes the paper.

4

## 2. Related Surveys

Other works have already surveyed contributions in malware classification. For example, Bazrafshan et al. [33] classify malware detection and classify methods into three types: signature-based, behaviour-based and heuristic-based methods. Also, they recognize five classes of features based on the proposed heuristic-based method: opcodes, API calls, control flow graphs, n-grams, and hybrid features. Another work presented by Shabtai et al. [34], which studies how to detect malware using static features. In this paper, we study more features (static and dynamic features) used for malware classification.

Ucci et al. [30] survey the literature on machine learning approaches for malware detection and analysis. They classify the surveyed articles into three categories: objectives (expected output), features, and algorithm used. They also highlight a set of problems and challenges and identify the new research directions. Similarly, the survey presented by [31] presents a comparative analysis on intelligence-based malware classification. In particular, they report cons, pros and problems associated with each machine learning-based malware classification technique. Souri and Hosseini [32] also provide a taxonomy of AI-driven malware detection techniques. Our paper looks at a larger range of articles by including many works on malware classification and composition analysis. We also include other works related to non-AI-driven classification techniques. Furthermore, We also include new challenges related to AI-driven malware classification techniques.

Also, Basu et al. [35] study different works relying on AI-powered malware classification techniques. In particular, they coin five types of features: a PI call graph, byte sequence, PE header and sections, assembly code frequency and system calls. Also, Ye et al. [36] study many different aspects of malware classification processes. More specifically, they spot the light on a number of issues such as incremental learning, and adversarial learning. Recently, Ori et al. [37] survey the literature on techniques used for dynamic malware analysis, which includes a description of each technique. In particular, they present an overview of machine-learning methods used to improve the capability of dynamic malware analysis. Compared to the above-motioned works, this paper determines the main issues and challenges on malware classification and composition analysis. Also, we identify a number of trends on the topic, with guidelines on how to improve solutions to address new and continuing challenges.

In addition, Barriga and Yoo [38] survey the literature on malware evasion techniques and their impact on malware analysis techniques. This paper extends beyond that and includes recent AI-driven works used to overcome malware evasion techniques.

## 3. Taxonomy of Malware Classification

We present in this section the taxonomy of malware classification. We define two categories (or dimensions) to organize the existing works. The first category presents the features that our work is based on. In particular, we discuss the different methodologies used for extracting features, e.g., dynamic and/or static techniques, and what types of features are used, e.g., assembly code. The second is concerned with the type of algorithm that is adopted for the detection and analysis, e.g., artificial inelegance-driven algorithm.

Figure 1 shows the proposed taxonomy. The rest of this section is organized as follows (according to the proposed taxonomy). Subsection 3.1 describes malware analysis features, while subsection 3.2 discusses existing algorithms.

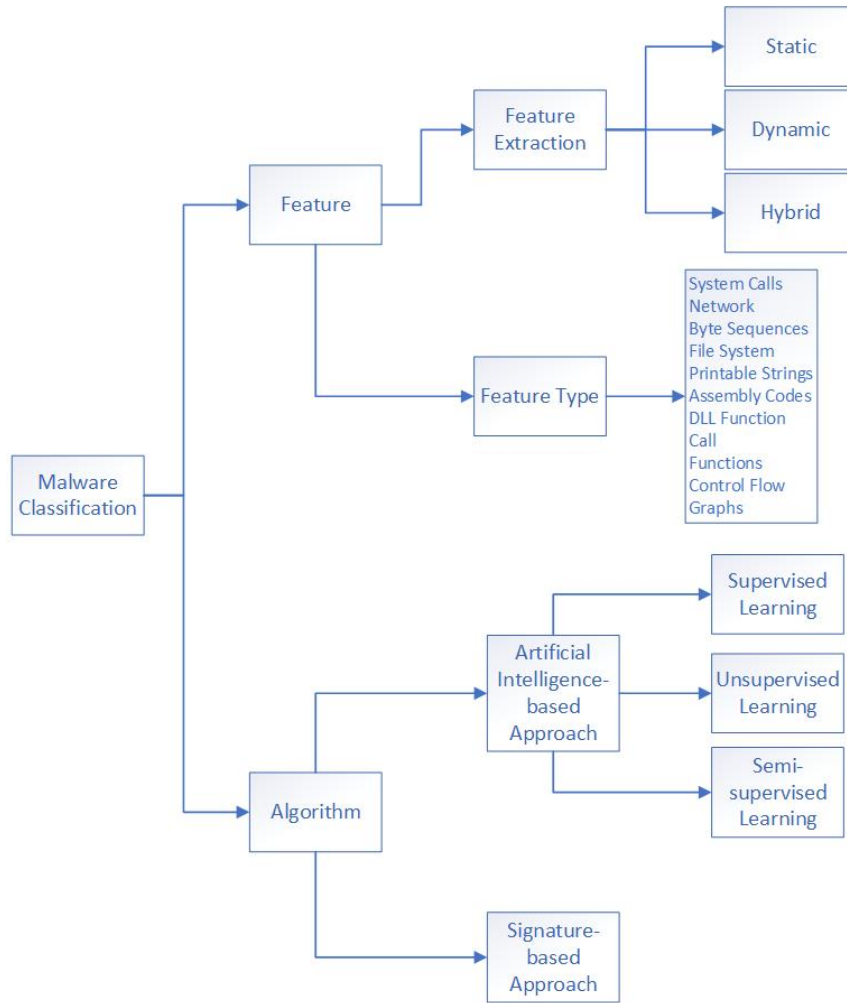Figure 1: The proposed taxonomy

## 3.1. Malware Analysis Features

This subsection presents the features of samples that are used for the analysis. In subsection 3.1.1, we show how features are extracted, while in subsection 3.1.2, we show type of features that are taken into account.

## 3.1.1. Feature Extraction Methods

In this section, we review the following three feature extraction methods: static, dynamic and hybrid methods.

*Static Method.* Static feature extraction is a method to extract features from the content of the executables without running them [39]. The static features can be extracted using the file format, e.g., Portable Executable (PE) and Common Object File Format (COFF) [12, 18, 22, 25]. The static features can also be extracted without any knowledge of the format. Features extracted this way can be byte sequences, file size, byte entropy, etc. [12, 17, 20, 25]. The advantage of the static feature extraction method is that it covers the complete binary content. But the problem is that static features are prone to packing and polymorphism since most of the features that are statically extracted come from encrypted contents rather than the original program body [40].

*Dynamic Method.* Dynamic feature extraction consists of running the executable usually in an insulated environment which can be a virtual machine (VM) or an emulator and then extract features from the memory image of the executable or from its behaviors [39]. Since malware equipped with packing and polymorphism has to exhibit the real malicious code to achieve their goals, dynamic feature extraction is more resistant to those malware techniques compared with static feature extraction method [40].

Anderson et al. [21, 41] use Xen [1] and Royal et al. [42], Dai et al. [19], and Islam et al. [22] use VMWare [2] to create their VMs and perform dynamic analysis. Kolosnjaji et al. [27] use Cuckoo sandbox [3] which is an open source automated malware analysis system to extract API calls. Other researchers who work for an anti-virus engine use the VMs as parts of their anti-virus engines to dynamically extract features [24, 26].

In fact, there are two categories of an emulator: a full-system emulator and application level emulator. A full-system emulator is a computer program that emulates every component of a computer, including its memory, processor, graphics card, hard disk, etc., with the purpose of running an unmodified operating system. Qemu [4] is a full-system emulator used by several systems [40, 43, 23]. Considering the time-consuming of full-system emulator, Cesare and Xiang [15] propose to use application level emulation to unpack malware more efficiently so that only the parts which are necessary

---

[1]https://www.xenproject.org/
[2]https://www.vmware.com/
[3]https://cuckoosandbox.org/
[4]https://www.qemu.org/

to execute the file including instruction set, API, virtual memory, thread and process management, and OS specific structures are implemented.

One problem of dynamic feature extraction methods is that it does not reveal all the possible execution paths [40]. Malware may have detection routines to check whether it is executed in a virtual machine or emulator. When malware finds itself executing in such an environment, it will halt its execution so dynamic models will fail to recognize it as malware. The methods to detect whether an executable is executed inside a VM can be found from several papers [44, 45]. Another problem of dynamic methods lies in its execution time which takes much more than static feature extraction [40].

*Hybrid Method.* This method is used to achieve higher detection rate by merging some of the static feature extraction characteristics with some of the dynamic feature extraction characteristics [39].

Our survey has revealed that most of the surveyed papers were based on the dynamic feature extraction approach [46, 47, 48, 49, 50, 51, 21, 52, 53, 54, 24, 55, 56, 57, 58, 59, 60, 61, 62, 63]. while the others adopt, in equal proportions, either the static approach alone [64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83] or a hybrid approach [84, 41, 85, 22, 23, 86, 47].

*3.1.2. Type of Features*

In this section, we classify the features that are used by malware analysts and explain how each type is practically extracted and represented.

*Printable Strings.* A printable string is a sequence of ASCII characters terminated with a null character. Schultz et al. [12] find that malware have some similar strings that distinguish it from and that Goodware also has some common strings that distinguish them from malware. Printable strings are represented as binary features, where "1" represents a string that is present in an executable and "0" represents that it is absent from all systems [12, 22, 24, 26].

Schultz et al. [12] extract printable strings from the headers of PE files. The extraction is straight-forward since the header is in plain text format.

Dahl et al. [24] and Huang and Stokes [26] extract null-terminated objects dumped from images of a file in memory [24, 26] as printable strings. The coverage of their methods is better than just extract printable strings from header [12] but their could be some false positive results.

Islam et al.[22] use the strings utility in IDA Pro [5] to extract printable strings from the whole file.

Different from other works, Saxe and Berlin [25] do not take printable strings as binary features but use their hash values and the logarithm of the string lengths to create a histogram and use the counts of printable strings in each bin of the histogram as features. They take all the byte sequences of length six or more that are in the ASCII code range as printable strings which is also slightly different from other works.

Essentially, the functionality of most malware does not rely on printable strings. Thus, when malware creators find that some strings accidentally are used by malware detectors, they can eliminate them or even if the printable strings are necessary, they can break them into characters that are distributed in different positions. Therefore, printable strings are not reliable features.

*Byte Sequences (Byte Code).* Executable files consist of byte sequences (also known as byte code). A byte sequence may belong to the metadata, code, or data of an executable file. As has been stated, byte sequences are important signatures of malware since malware may share some common sequences that are exactly the same or follow the same regular expression. Thus, byte sequences are also appropriate to be features for malware analysis systems [12, 17, 41, 25].

Schultz et al. [12] use bigram byte sequences in the form of binary features and they claim byte sequence feature is the most informative feature because it represents the machine code in an executable. In fact, this is not entirely true since some byte sequences come from metadata or data section. Even if a byte sequence is from code section, since instructions have variable length in some architectures, byte sequences may not match machine code. And their byte sequence feature has the problem of dimension explosion since there are too many different bigram byte sequences and it is too large to fit into memory so they could only split the byte sequence set into several sets and feed them to multiple native bayes models.

To solve the dimension explosion problem, Kolter and Maloof [17] use information gain to select the top 500 informative 4-gram byte sequences as binary features from 255 million distinct 4-grams.

Different from the above two works, Anderson et al. [41] do not use byte sequences per se as features but fit byte sequences into a Markov Model so

---

[5]https://www.hex-rays.com/products/ida/

essentially the feature they use is transition probability from one byte to another.

Chen et al. [25] use the byte entropy of each 1024 byte window and the occurrence of each byte to form a histogram and evenly separate each axis into 16 bins to form a 256 length feature vector.

Nataraj et al. [20] convert the whole byte sequence of a file into a picture in which each byte represents the grey scale of a pixel. They find that the malware that belongs to the same family appear very similar in layout and image. The width of the image that is used to transform the 1D byte sequence into a 2D matrix is determined by the size of the file. The image feature of the malware image is computed using the algorithm proposed by Oliva and Torralbat [87]. The main advantage of image-based techniques is that they are robust against many types of obfuscations [88].

Byte sequences are not reliable in most cases. This is due to the fact that obfuscation techniques such as instruction substitution and register reassignment can change the opcodes and oprands respectively, which means that the machine code is changed. In all these works, the byte code is statically extracted but the main program body encrypted with different algorithms or keys through Packing and Polymorphism will change the byte sequences.

*Assembly Code.* Machine code and assembly code can be translated to one another through assembly and disassembly. Assembly code has some advantages over machine code as a feature for malware analysis. First, assembly code can be understood by a programmer and therefore as a kind of feature, assembly code is more convenient to be preprocessed (e.g., grouped into categories according to the function, filtered, truncated etc.) to appear as a more informative feature. In addition, malicious code is often encrypted by packing or polymorphism so it is impossible to get it from the original byte sequence, however, dynamically extracted assembly code has been decrypted so it includes the malicious code.

Moskovitch et al. [18] propose that assembly code can be more robust than machine code for the analysis of malware since the same malicious engine may locate in different locations of a file, and thus may be linked to different addresses in RAM or even perturbed slightly so by dropping the oprands and just using opcode the robustness is improved. They extract assembly code by dissembling the executables with IDA Pro. They try both term frequency (TF) and term frequency–inverse document frequency (TF-IDF) of each opcode n-gram (n=1,2,...,6) as features and use document fre-

quency (DF), information gain ratio, or Fisher score to select features. Their best result is achieved using TF values of opcode bigram as features filtered by Fisher score. One disadvantage of their method is that it is still prone to dead code insertion, operation transpositions, packing, and polymorphism. Another one is dropping operands causes loss of information which may subsequently lead to loss of precision.

To counter packing and polymorphism, Dai et al. [19] run malware in a VM and record the sequence of the running byte code which will be disassembled to assembly code. They use three kinds of two-opcode combinations: unordered opcodes in a block, ordered but not necessarily consecutive opcodes in a block, consecutive opcodes in a block. This way their features is more resistant to dead code insertion and reorder of operations. They use the association between the frequency of a feature in training dataset and a class as criterion and apply a variant of Apriori [89] to select top $L$ features. Although unordered opcodes and ordered (but not necessarily consecutive opcodes) in a block improve the resistance to dead code insertion and reorder of operations, those features are too flexible so they also bring more false positive situations.

Royal et al. [42] is another work aiming to detect code that is hidden and can only be seen dynamically. The way they do it is to store the static code of an executable and check whether each operation executed is within the stored static code area. If it is not, it is a part of hidden-code. They claim that the main malware engine should be in the hidden-code if both of them exist and experiment results also illustrate the hidden-code enhances the accuracy of ClamAV [6] and McAfee Antivirus [7].

Anderson et al. [21, 41] use the transition probability from one opcode to another as features, which is similar to how they use byte sequence feature. In their paper [21], they just extract assembly code by recording the execution of an executable in a VM which is similar to the way Royal et al. [42] use. In their second paper [21], they also use IDA Pro to disassemble the executable, and the assembly code from the two sources are used as two independent feature sets. In addition, they also group instructions into categories in several granularities according to the functions of the instructions to reduce the impact of instruction substitution in their second paper [21]. In their

---

[6]http://www.clamav.net/
[7]https://www.mcafee.com/en-us/index.html

12

preliminary experiment, they also find if they use instructions with oprands, the performance will be worse [21].

Santos et al. [23] disassemble executables to acquire their assemble code and then use weighted opcode n-gram frequencies as one of their features. The weight is the product of the information gain of all opcodes in the n-gram times the normalized TF of the n-gram.

*API/DLL System Call.* DLL files and functions of DLL files used by an executable expose the system services they use. Native system calls and Windows API calls an executable invokes are shown by the functions of DLL files it depends on. Therefore, what behaviors it may intend to do or what it would be able to do can be inferred.

Schultz et al. [12] extract the DLL files by an executable used, the functions in DLL files, and the number of function of each DLL as features from metadata in order to understand how resources affected an executable's behavior and how heavily each DLL is used. The first two are used as binary features and the third is a real-valued feature.

Bayer et al. [40] and Santos et al. [23] extract calls to Windows API functions dynamically using an emulator. Then, they use those API functions to acquire actions of an executable during execution including I/O activity, registry modification activity, process creation/termination activity, network connection activity of an executable, self-protection behavior, system information stealing, errors caused by the execution, and interactions with Windows Service Manager.

Fredrikson et al. [43] also use an emulator to monitor system calls. Then, they use the relations between system calls and their parameters to form a dependency graph in which nodes are system calls and edges connect system calls sharing some parameter. They define a behavior to be a subgraph of it and behaviors that can be adopted to distinguish malware from Goodware will be mined and used to detect malware.

Anderson et al. [41] and Huang and Stokes [26] group the system calls into high-level categories where each category represents functionally similar groups of system calls, such as painting to the screen or writing to files. Anderson et al. [41] then feed the trace of groups of system calls to a Markov chain so that they use transition probability of system calls to be the feature. Huang and Stokes [26] use those high-level API call events as binary features.

Islam et al. [22] and Dahl et al. [24] extract Windows API function calls and their parameters by running an executable in a VM. Islam et al. [22]

treat Windows API functions and parameters as separate entities and use the occurrence frequency of each entity as their feature. Dahl et al. [24] use combination of a single system API call, one input parameter, and API tri-grams which consist of three consecutive API function calls, as binary features which are subsequently selected using mutual information.

Kolosnjaji et al. [27] use the dynamic malware analysis system Cuckoo sandbox to extract the sequence of the Windows system calls invoked by an executable. They use one-hot representation of them and feed the full sequence of system calls with the order to a sequential deep learning model.

Similar to assembly code, Windows API call sequences can also be obfuscated. For instance, malware authors can make an executable invoke some irrelevant API calls and submerge the API calls they use to fulfill their purpose in them. Thus, this feature is not reliable in most cases.

*Control Flow Graphs.* A control flow graph is a directed graph that represents the flow of the program, where nodes are the instructions while the edge between two nodes represents the order of sequence of execution of the two instructions. A vertex in the graph is a basic block in the middle of which there is no jump or branch instructions. A directed edge represents jumps in the control flow. Control flow graphs are used as features or signatures to detect malware in several papers [15, 41].

Cesare and Xiang [15] state that similar malware usually have similar high-level structured control flows. They find that compressed and encrypted data have relatively high entropy so they first use entropy of byte sequence to detect whether an executable is packed or not. If so, they use an application level emulator to extract hidden code. They still use entropy of byte sequence to detect completion of hidden code extraction. Then the memory image of the binary is disassembled using speculative disassembly [90]. Finally, they use the process of structuring to recover high-level structured control flows from control flow graphs of procedures and represent them using strings of character tokens. The strings representing control flow graphs are all saved as signatures. An example of the relation between a control flow graph and the signature string is shown in Figure 2.
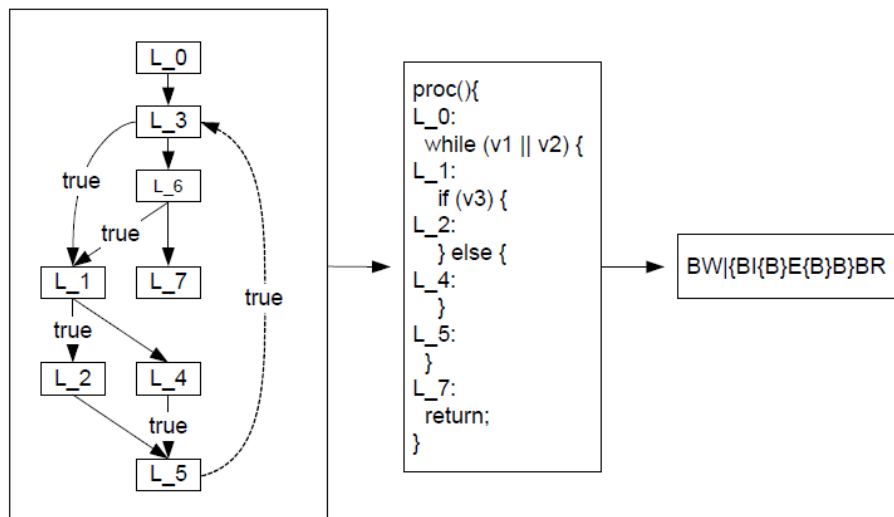
Figure 2: The relationship between a control flow graph, a high level structured graph, and a signature.

Anderson et al. [41] also find that it is largely not easy for a polymorphic virus to build a semantically similar version of itself while changing its control flow graph enough to avoid detection. Therefore, they use control flow graphs as features. More specifically, they use the occurrence frequency of each k-graphlet (a subgraph of k nodes) in the control flow graph to represent control flow graph.

To counter the detection using control flow graphs, malware authors can use control flow flattening and bogus control flow obfuscation techniques to change the control flow without affecting the functionality so that the effectiveness of control flow graph feature will be harmed [91, 92].

*Function.* Some papers (e.g., Islam et al. [22] and Chen et al. [14]) use function level features for malware classification.

In particular, Islam et al. [22] find function length that consists of statistically useful information in distinguishing between families of malware. After obtaining the assembly code of each executable, they calculate the length of them by measuring the number of bytes of code and use the occurrence frequency of each function lengths as a feature. However, obviously, function length is the least robust feature against obfuscation. Function length can be arbitrarily increased by inserting dead code or decreased by splitting them

15

into multiple functions.

One should note that two functions which are semantically similar to each other are considered to be clones of each other. To this end, Chen et al. [14] assume that some files that belong to the same malware family share some functions which are connected using clone relation. So they cluster functions to groups in which any two functions can be connected directly or indirectly using clone relation and pick one function from each group as an exemplar to be a signature. They use NiCad [93] to detect whether two functions are clone to each other. However, to use one function to represent a group of functions is problematic. Since the same function evolves over generations, the newest version may look quite different from the original one. If the older version is picked as the exemplar, the clone detector may fail to identify some unknown new generation of it. Although their system works on Android APK files, the methodology can be directly applied to classifying executable malware.

*Miscellaneous File Information.* Some miscellaneous file properties can help engineers distinguish malware from Goodware since the average or majority values of them are significantly different between the two groups. So that those properties are also used as features. They are file size [40, 41], exit code [40], time consumption [40], entropy [94][41], packed or not [41], number of static/dynamic instructions [41], and number of vertices/edges in control flow graph [41]. These features may be helpful but obviously not informative enough.

*Conclusive Remarks.* The effectiveness of using all the aforementioned features can be somehow diminished or they are not informative enough. So many papers use multiple features [12, 41, 24, 22, 23, 25, 26]. The intuition is that any single feature source can be obfuscated to evade the detection but it is extremely difficult to obfuscate all features simultaneously without hindering the functionality [41, 22].

### 3.2. Malware Classification Algorithms

The extracted features introduced in the previous section are fed into malware detection/classification systems. They can be categorized as signature-based approaches and artificial intelligence-based approaches.

### 3.2.1. Signature-based Approaches

Signature-based detection is the most papular approach used in most antivirus engines. Those signatures are created by human malware defenders

16

through examining the collected malware samples [12, 13]. More specifically, the antivirus engines detect or classify malware by checking whether the files to be analyzed contain malware signatures. The signatures of malware can take many formate including filename, text strings, or regular expressions of byte code [12, 13]. Signatures are usually also hashing of the entire file. One should note that signature-based techniques can only detect malware originates from known malware which does not change significantly. As a result, attackers can exploit these techniques by hiding the malicious behaviour of malware using anti-analysis techniques such as packing, obfuscation, polymorphism, and metamorphism (Section 8.1 provides more details about these techniques). Therefore, the code looks quite different from its original version. The main shortcoming of signature-based method is it has high precision but low recall and the other one is labor-intensive.

Some works [15, 14, 16, 15, 14] address the problem of manual signature crafting by proposing automatic signature generation techniques. The content of the signatures can be windows system call combinations, control flow graph, and functions.

### 3.2.2. Artificial Intelligence-based Approaches

The section discusses artificial intelligence-based malware classification approaches. These approaches can be categorized as traditional machine learning models, deep learning models, association mining, graph mining and concept analysis, and signature creation and search methods. The existing artificial intelligence-based approaches also can be classified according to the learning method used as follows: supervised, unsupervised or semi-supervised.

In a supervised malware classification model [64, 50, 58, 60, 81, 63, 59, 46, 95, 96, 65, 22, 55, 23, 74, 80, 62, 82, 71, 97, 85, 72, 24, 25, 67, 54, 76, 98, 99, 21, 69, 57, 61], the classification algorithm learns on a labeled dataset, which enable the algorithm to evaluate its accuracy on training data. In contrast, an unsupervised malware classification model [75, 83, 49, 62, 100, 101, 102, 47, 84, 53, 69], provides unlabeled data that the algorithm tries to make sense of by extracting patterns without guidance. Semi-supervised malware classification models [68, 103, 75, 78] combine both labeled and unlabeled data.

*Traditional Machine Learning Models.* The most popular traditional machine learning models used by surveyed papers are Naive Bayes classifier (NBC)

[64, 65, 50, 58, 60, 81, 63], rule-based classifier[64, 59, 81, 46, 95, 96], decision tree (DT) [65, 96, 50, 22, 55, 72, 23, 74, 58, 60, 80, 62, 82], K-nearest neighbors (K-NN)[71, 62, 97, 96, 50, 22, 60, 72], Bayesian Network [85, 72, 23], Neural Network (NN) [24, 25], Random Forest (RF) [67, 54, 22, 58, 76, 60, 80, 98, 99, 63], Hidden Markov Models (HMM) [104, 105, 106, 9] and Support Vector Machine (SVM) [65, 96, 50, 21, 69, 54, 22, 71, 72, 23, 57, 58, 76, 60, 61, 62, 81, 63]. Those papers which use traditional machine learning models normally try multiple machine learning models [12, 17, 18, 19, 22, 23].

Below, we briefly introduce the above mentioned machine learning models.

**Naive Bayes Classifier (NBC)**   An NBC [107] uses Bayes' theorem to determine the conditional probability of a sample belonging to a class given the input features which can be formally described in the following equation:

$$P(C_i|x) = \frac{P(x|C_i)}{P(x)}P(C_i) \tag{1}$$

where $x$ is a sample and $C_i$ is the probability the sample belongs to class $i$. It is based on the Naive Bayes conditional independence assumption that all the features are independent to each other given the class it belongs to:

$$P((x_1, x_2, ..., x_n)|C_j) = P(x_1|C_j)P(x_2|C_j)...P(x_n|C_j) \tag{2}$$

where $x_j$ is a feature of $x$. Although the assumption do not hold, the prediction results are good in many occasions and the result is explainable which means how much each feature contributes is visible.

**Decision Tree (DT)**   A DT classifier [108] uses a tree structure to represent the classification process. Internal nodes of a DT are tested on the values of features and edges correspond to a choice on values of a variable. Leaf nodes represent the final class of samples fall into it. The tree structure is constructed based on the informativeness of each feature conditioned on the current choices such as information gain ratio and Gini index. A DT is also an interpretable classifier and a DT can be translated sets of if-else-then rules.

**K-Nearest Neighbor (KNN)**   A KNN [109] is an instance-based classifier. The model finds the K nearest neighbors of a given sample with some

18

distance metrics (e.g., Euclidian, cosine), and predict it to be the (weighted) majority vote of the classes of the k nearest neighbors.

**Support Vector Machine (SVM)**  An SVM [110] is a binary classifier which calculates a hyperplane that separates samples from two classes with the largest margin. An important characteristic of an SVM is it can utilize kernel trick to map samples from the original feature space to a high-dimensional (even infinite) feature space to perform non-linear classification.

**Bayesian Network (BN)**  A BN [111] is a probabilistic graphical model which represents variables as vertices and the dependencies as directed edges. The graph is used for the inference of probability of any variable.

**Rule-based Classifier**  A rule-based classification [112] refers to any classification method that allows us to use of IF-THEN rules for prediction. An example of a rule-based classification is RIPPER [113], which is used to build a set of rules to classify samples while minimizing the error of the number of misclassified training samples.

**Neural Network (NN)**  An NN [114] is a biologically-inspired programming paradigm that allows a computer to learn from observational data. It consists of a network of functions (i.e., parameters) which enables the computer to learn, and to fine tune itself, through analyzing new data.

**Random Forest(RF)**  An RF classifier [115] constructs a set of DTs from the subset of training set (selected randomly). The votes are then aggregated from trees in order to decide the final class of the test sample.

*Deep Learning Models.* Deep learning models allow us to automatically abstract and extract robust and useful features for efficient and reliable malware classification. This can be done using multiple layers of abstraction to learn the "good" representation of the data [116]. An example of deep learning models are autoencoder [117], stacked denosing autoencoder [116], restricted Boltzman Machine (RBM) [118].

Dahl et al. [24] applies their 179,000 binary features to a deep learning model. The first layer is a random projection layer which maps the input features to a much lower dimensional space (4000 dimension). The difference

between the random projection layer and a normal fully connected layer is the weight of the projection matrix is not updated. The entries of it are sampled following an independent and identically distribution over -1,0,1. On top of that, they apply 1 to 3 fully connected layers with sigmoid activation functions and a 136-way softmax layer as output. They also try using a Gaussian-Bernoulli restricted Boltzmann machine (RBM) to pre-train the hidden layers. The best result is achieved by the model with 1-hidden layer without pre-training which is 9.53% test error rate. They also find the random projection performs better than Principal Component Analysis (PCA).

Saxe and Berlin [25] propose a deep feed-forward neural network consisting of four fully connected layers, where the dimensions of the first three layers are 1024 followed by a dense layer to get the output. They apply dropout to the first three layers. The activation functions of the first two layers are parametric rectified linear units (PReLU) to yield improved convergence rate without loss of performance and the activation function of the third layer is sigmoid. They also use Bayesian Calibration to calculate the unbiased probability that an executable is malware. They achieve a detection rate of 95% and a false positive rate of 0.1% on a dataset of 431,926 samples.

Huang and Stokes [26] propose a neural network for multi-task training. One task is a malware detection to predict whether an unknown software is malicious or benign and the other is to predict if it belongs to one of 98 important malware families. Huang and Stokes [26] also use a random projection layer to reduce the dimension to 4,000 from 50,000 and then they normalize each of the 4,000 dimension to be zero mean and unit variance. Then they use 4 hidden layers with dropout and RELU activation. On top of it is two single layers for each of the two classification task. The final loss function is a weighted sum of each of the individual loss functions. Experiment results show that multi-task learning only improve the performance of malware detection and harm the performance of malware classification in most experiment settings. Specifically, the best result for malware detection is 0.3577% test error which uses two hidden layers and multi-task learning and the best result for malware classification is 2.935% test error which uses one hidden layer and either single task or multi-task learning.

Kolosnjaji et al. [27] propose a combination of convolutional neural network (CNN) and Long Short-Term Memory (LSTM) networks to predict the family of an executable using the dynamically extracted system call sequence. They first use two convolution layers to capture the correlation between consecutive API calls and then apply max-pooling to reduce the dimensionality.

The output sequence is fed to a LSTM layer to model the sequential dependencies of API calls. Then a mean-pooling layer is used to extract important features from the LSTM output. They also use Dropout to prevent overfitting and a softmax layer to output the probability of each class. Their proposed deep learning model significantly outperforms feed-forward neural networks, CNN, SVM, and Hidden Markov Model and achieves 85.6% on precision and 89.4% on recall. The advantage of their model is it can fully utilize the order of system calls which may also be a drawback if the system call sequence is obfuscated. One problem of their model is they use mean-pooling rather than max-pooling to extract features of highest importance produced by LSTM is not quite reasonable.

*Associative Classifier.* An associative classifier relies on association rules that can be used to distinguish samples between two classes to perform classification. It is a special case of association rule mining where only the class of a sample can be the consequent (a.k.a. right-hand-side) of a rule. Ye et al. [16] proposes to use hierarchical associative classifiers (HAC) to classify executables based on API calls. There are three techniques regarding the creation of an associative classifier: 1) adopt FP-Growth algorithm to find candidate association rules (i.e., combination of API calls) 2) prune the candidate rules based on $\chi^2$, data coverage, pessimistic error estimation, significance w.r.t to its ancestors 3) reorder rules: first rank the rules whose confidences are 100 by confidence support size of antecedent (CSA) and then re-order the remaining rules by $\chi^2$ measure. Using those three techniques, they create a 2-level associative classifier to detect malware from a gray list labeled by a signature-based anti-virus engine. The first-level associative classifier is aimed for higher recall of malware. It only keeps the rules of Goodware with 100% confidence and the rules of malware with confidence greater than a pre-defined threshold; then it uses the rule pruning technique to decrease the generated rules and create the classifier; finally uses "Best First Rule" technique to find samples from the gray list. The samples labeled to be malware by the first associative classifier are fed to the second level associative classifier which is aimed at optimizing the precision. It works with the following steps: select those samples whose prediction rules of malware have 100% confidences, marking them as "confident" malware; ranking the remaining minority class files in an descending order based on their prediction rules' $\chi^2$ values; select the first k files from the remaining ranking list and marking them as "candidate" malware; mark the remaining files as "deep gray" files.

Experiment results show the proposed HAC is effective. In addition, HAC is also an interpretable classifier which can be easily represented as simple if-then rules.

*Graph Mining and Concept Analysis.* Fredrikson et al. [43] extract behaviors (dependency graphs of system calls and their parameters) that can distinguish malware from Goodware using structural leap mining [119]. Then they use the behaviors to form discriminative specifications. A specification is a set of behaviors and a characteristic function that describes one or more subsets of the set. A software matches a specification if it matches all of the behaviors in at least one characteristic subset. A specification is entirely discriminative if it matches malicious software but does not match benign software. They use formal concept analysis [120] and Simulated Annealing algorithm [121] to find an approximate optimal specification which has true positive larger than a threshold and lowest false positive among all specification larger than that true positive rate. During test, if a program matches a specification, it will be classified to be malware. The created specification can be used in the detection of unseen malware with a 86% true positive rate and 0 false positives on a dataset of 961 samples.

*Signature Search Methods.* Cesare and Xiang [15] first convert the control flow graphs of each procedure in an unkown executable to character strings in the same way they create signatures. Each procedure is assigned a weight using the length of its string:

$$weight_x = \frac{\text{len}(s_x)}{\sum_i \text{len}(s_i)} \tag{3}$$

Then they use BK Trees to retrieve the strings in the signature database which have less Levenshtein distance with strings representing procedures of the target file than a threshold. For a particular malware, once a matching graph is found, this graph is ignored for subsequent searches of the remaining graphs in the input binary. If a graph has multiple matches in a particular malware and it is uncertain which procedure should be selected as a match, the greedy solution is taken. The graph that is weighted the most is selected. For each malware that has matching signatures, the similarity ratios of those signatures:

$$w_{ed} = 1 - \frac{ed(x,y)}{\max(\text{len}(x), \text{len}(y))} \tag{4}$$

22

are accumulated proportional to the weights of the procedure. The final similarity between the unknown executable and a malware in the database is the product of two asymmetric similarities: a similarity that identifies how much of the input binary is approximately found in the database malware, and a similarity to show how much of the database malware is approximately found in the input binary. If the program similarity of the examined program to any malware in the database equals or exceeds a threshold of 0.6, then it is deemed to be a variant. Experiment results show that their method achieves 86% detection rate with 0 false positives which is better than 55 for commercial signature-based antivirus (AV) and 62-64 for behavior-based AV. Since they use a symmetric similarity calculated as the product of two asymmetric similarities, it can not handle asymmetric situations. For instance, if a very large unknown executable contains the whole program of a malware sample in the database but that malicious program only take up 1% of its whole content, the similarity would still be small and it can not be predicted to be malware.

Chen et al. [14] uses NiCad [93] to detect whether an APK file contains any function that is clone of an exemplar function which represents a signature of a malware family. If a match is found, the file is predicted to be an instance of that malware family. They achieve 96.88% accuracy on a dataset of 1170 APK files from 19 malware families.

## 4. Taxonomy of Composition Analysis Techniques

This section introduces the taxonomy of malware composition analysis techniques. We identify two major dimensions along which surveyed papers can be conveniently organized. The first one shows the steps used for composition analysis. The second dimension identifies the objective (i.e., strategy) of the analysis. Figure 3 shows a graphical representation of the proposed taxonomy.
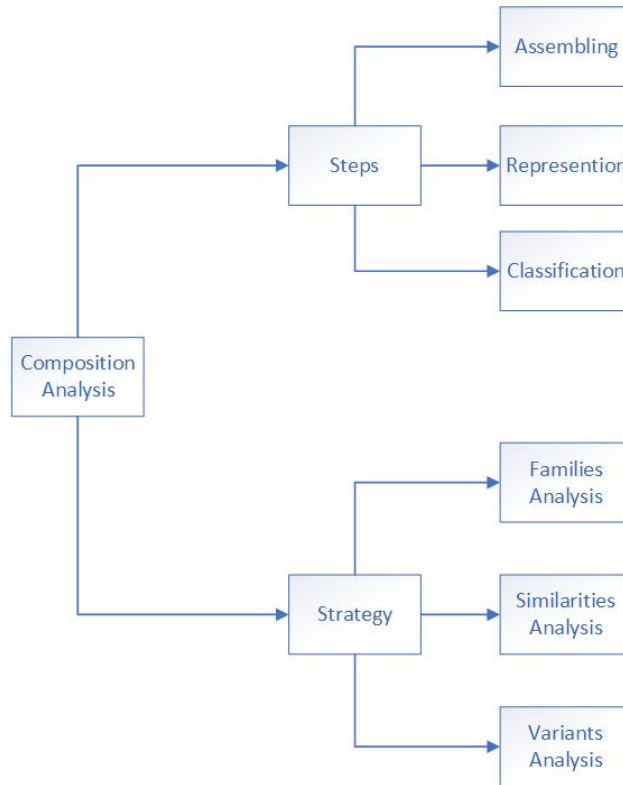
23

Figure 3: The proposed taxonomy

## 4.1. Steps

Composition analysis allows reverse engineers to analyze the composition of malware samples in order to understand their functionalities and behaviours. This, in turn, allows engineers to discern the intent of malware samples and the attackers. Moreover, it allows reverse engineers to rank the malware by severity and allows them to effectively triage their resources.

Basically, there are three main steps used for composition analysis: disassembling, representation, and classification.

### 4.1.1. Disassembling

Most software programs are delivered to users with compiled executables, rather than source code. Disassemblers make it feasible for reverse engineers to analyze software programs without source code. Technically speaking, a disassembler is a process of converting or translating machine language into

24

assembly language. The inverse operation of "disassembler" is an "assembler". There are many tools used for this purpose (e.g., IDA Pr [8]).

Disassembly methods can be categorized into the following two classes: static techniques and dynamic techniques. Methods that belong to the first class analyze the binary components statistically, parsing the opcodes in the binary file. Methods belong to the second class monitor the execution traces of a program in order to identify the instructions and recover disassembled version of the binary.

Both dynamic and static methods have pros and cons. Static analysis takes into consideration the whole program, while dynamic analysis can only focus on the executed instructions. As a result, it is not easy to ensure that the entire executable was visited when adapting dynamic analysis. However, dynamic analysis guarantees that the output (i.e., disassembly output) only contains actual instructions.

Generally speaking, there are two approaches for static analysis techniques. The first approach is called linear sweep [122]. This approach begins at the first byte of the binary and starts decoding one instruction after another. The main shortcoming of using liear sweep disassemblers is the high probability of errors which result from data embedded in the program. The second approach is called recursive traversal [123], which allows engineers to fix the problem of "embedd data" by following the Control Flow (CF) of the program [15, 41]. However, the problem with this approach is that it could fail to successfully analyze parts (i.e., functions) of the code. This is due to the fact that a control transfer instruction (e.g., jump) cannot be determined statically. This problem can be addresses by using a linear sweep algorithm to analyze unreachable regions in the code [124].

### 4.1.2. Representation Learning

The success of any malware classification and composition analysis technique generally depends on data representation. Although specific domain knowledge may help engineers design representations and a feature vector for an executable, a manual feature engineering process fail to consider the relationships between features and define those unique patterns that can distinguish executables.

Indeed, representation learning is a set of methods and/or techniques that

---

[8]https://www.hex-rays.com/products/ida/

25

enables a system to automatically extract the representation needed for malware classification from raw data (i.e., assembly code). This process replaces manual feature engineering and enables a malware classification system to learn the useful features and integrates them to perform a classification.

The motivation behind using feature learning is the fact that composition analysis methods often need inputs that are robust against anti-analysis techniques such as obfuscation and packing.

Deep learning approaches (e.g., stacked autoencoders [125], stacked Denoising autoencoders [116], Deep belief networks [126], ...) are known and considered as the (best) approaches for extracting robust features, which are used for building robust malware and similarity analysis tools for large-scale heterogeneous environment.

### 4.1.3. Classification

After disassembling executable samples, the assembly code functions are used to feed a representation learning module in order to obtain robust features and "good" representation of data. The function representation are then fed into any classification algorithms such as Naive Bayes classifier (NBC) [64], rule-based classifier[64], decision tree (DT) [65], K-nearest neighbors (K-NN)[71], Bayesian Network [85], Neural Network (NN) [24], Random Forest (RF) [67], Hidden Markov models (HMM) [127], and Support Vector Machine (SVM)[65]. The classification method enables us to identify the relationships between functions taking into account the following three analysis strategies: variants analysis, similarities analysis, and families analysis.

*Variants Analysis (VA).* VA [79, 59, 80, 83, 46, 47] enables engineers to realize that a malware sample is actually a variant of a known malware in the repository. This strategy allows us to understand to which extent malware have been evolved over time.

*Similarity Analysis (SA).* SA [48, 53, 49, 56, 128] allows engineers to recognize what parts (i.e., functions) of a malware sample are similar to known functions in the repository. This strategy allows us to focus only on new parts and prevent unnecessary investigation.

*Families Analysis (FA).* FA [101, 51, 102, 24, 70, 22, 71, 55, 76, 60, 61, 62, 97]. enables engineers to associate undefined malware to defined families. This strategy works under the assumption that malware from the same family

are similar to each other in terms of functionality. The difficulty to recognize them comes from the fact that some malware authors use anti-analysis techniques (e.g., obfuscation, packing, polymorphism, and metamorphism) to conceal that similarity.

## 5. Characterization of Surveyed Papers

In this section, we characterize each reviewed paper. Table 1 provides information about both algorithms and features used for each paper and highlights the main limitations. The table also shows the scalability of each work in terms of its ability to work in the presence of incremental update of the repository. The last column shows whether the proposed classification techniques are robust against anti-analysis techniques or not. As can be seen in the Table 1, most of the works use more than one classification algorithm for detecting and classifying malware in order to guarantee more accurate results. In Table 2, different approaches are compared w.r.t the of the main objective: malware detection and similarity analysis, families analysis and variants analysis.

Table 1: Summary of Extraction Methods, Classification Methods, and Limitation in Malware Classification.

| Begin of Table | | | | | |
|---|---|---|---|---|---|
| Work | Classification method | Features | Limitations | Scalability (Yes/No) | Robust against noisy inputs (Yes/No) |
| [129] | k-NN and SVM | Byte Code | Not robust against unseen inputs | Yes | No |
| [130] | NN | Byte Code | Vulnerable to adversarial attacks | Yes | Yes |
| [131] | k-NN and NN | Byte Code | Vulnerable to adversarial attacks | Yes | Yes |

| | | | | | Continuation of Table 1 | | | |
|---|---|---|---|---|
| **Work** | **Classification method** | **Features** | **Limitations** | **Scalability (Yes/No)** | **Robust against noisy inputs (Yes/No)** |
|---|---|---|---|---|---|
| [65] | DT, Naïve Bayes, and SVM | Byte Code | Not robust against noisy inputs | Yes | No |
| [132] | k-NN, NN, and SVM | Byte Code | Vulnerable to adversarial attacks | Yes | Yes |
| [73] | RF | Miscellaneous File Information | Needs a large number of labeled examples (malicious and benign) | Yes | Yes |
| [74] | DT, RF | Miscellaneous File Information | Works only under the assumption that the new samples are not packed | Yes | No |
| [57] | SVM | Internet Traffic | Not scalable (tested using vary small datasets) | No | Yes |
| [75] | Cluster Analysis | Miscellaneous File Information | Unable to classify new examples/samples | Yes | No |
| [64] | NBC | Printable Strings and Byte Code | Not robust against noisy inputs | Yes | No |
| [96] | DT, NBC, SVM | API | Not scalable (tested using very small datasets) | No | Yes |

| Continuation of Table 1 | | | | | |
|---|---|---|---|---|---|
| **Work** | **Classification method** | **Features** | **Limitations** | **Scalability (Yes/No)** | **Robust against noisy inputs (Yes/No)** |
| [103] | BN | Miscellaneous File Information | Not efficient giving new samples | Yes | No |
| [50] | DT, NBC, SVM, k-NN, NN and SVM | API and Miscellaneous File Information | Not scalable (tested using small datasets) | No | Yes |
| [21] | SVM | Byte Code and API | Not scalable (tested using very small datasets) | No | Yes |
| [41] | SVM | Byte Code, Assembly Codes and API | not scalable (tested using very small datasets) | No | Yes |
| [85] | BN | API | Not robust against noisy inputs | Yes | No |
| [23] | BN, DT, k-NN classification, SVM | Assembly Codes and API | Not robust against noisy inputs | Yes | No |
| [58] | DT, RF, Naïve Bayes,SVM | Byte Code and API | Not scalable (tested using very small datasets) | No | Yes |
| [78] | BN | Miscellaneous File Information | Not robust against unseen inputs | Yes | No |
| [59] | Rule-based classifier | API | Not scalable (tested using very small datasets) | No | Yes |

| | | | | | |
|---|---|---|---|---|---|
| | | Continuation of Table 1 | | | |
| **Work** | **Classification method** | **Features** | **Limitations** | **Scalability (Yes/No)** | **Robust against noisy inputs (Yes/No)** |
| [98] | RF | Internet Trafic | Not robust against unseen inputs | Yes | No |
| [99] | RF | API and Miscellaneous File Information | Not robust against noisy inputs | Yes | No |
| [25] | NN | Printable Strings and Miscellaneous File Information | Not robust against noisy inputs and not scalable (tested using very small datasets) | No | yes |
| [46] | Rule based classification | API and Miscellaneous File Information | not scalable (tested using very small datasets) | No | Yes |
| [47] | Cluster analysis | API and Miscellaneous File Information | Requiring user interactions | Yes | No |
| [101] | Cluster analysis | Byte Code | Not scalable (tested using small datasets) | No | Yes |
| [51] | Matching (graph theory) | API | Not robust against noisy inputs | Yes | No |
| [102] | Cluster analysis | Assembly Codes | Not robust against noisy inputs | Yes | No |
| [24] | NN | Byte Code and API | High error rate | Yes | No |

| Continuation of Table 1 | | | | | |
|---|---|---|---|---|---|
| **Work** | **Classification method** | **Features** | **Limitations** | **Scalability (Yes/No)** | **Robust against noisy inputs (Yes/No)** |
| [70] | Clustering | Assembly Codes | Not robust against noisy inputs | Yes | No |
| [22] | DT, k-NN classification, RF, SVM | Byte Code and API | Not robust against unseen inputs | Yes | No |
| [71] | k-NN classification and SVM | Assembly Codes and Miscellaneous File Information | Not robust against unseen inputs | Yes | No |
| [55] | DT | Internet Traffic | Not scalable (tested using very small datasets) | No | Yes |
| [76] | SVM, RF and DT | Internet Traffic and Byte Code, Assembly Codes and API | Not robust against noisy inputs | Yes | No |
| [61] | SVM, RF and DT | Internet Traffic and Byte Code and API | Not scalable (tested using very small datasets) | No | Yes |
| [60] | DT, RF, k-NN classification and NBC | API | Not robust against unseen inputs | Yes | No |
| [62] | DT, k-NN classification and SVM | Miscellaneous File Information and network | Not robust against noisy inputs | Yes | No |
| [133] | k-Means | Assembly Codes | Not robust against noisy inputs | Yes | No |

| Continuation of Table 1 | | | | | |
|---|---|---|---|---|---|
| **Work** | **Classification method** | **Features** | **Limitations** | **Scalability (Yes/No)** | **Robust against noisy inputs (Yes/No)** |
| [48] | Hierarchical Clustering | API, Miscellaneous File Information, and Internet Traffic | Not scalable (tested using very small datasets). Not robust against noisy inputs | Yes | No |
| [49] | Cluster analysis | API | Not robust against noisy inputs | Yes | No |
| [53] | Cluster analysis | Byte Code and API | Not robust against noisy inputs | Yes | No |
| [56] | NN | API | Not robust against noisy inputs and not scalable (tested using small datasets) | No | Yes |
| [72] | DT, k-NN classification, BN and RF | Assembly codes | not scalable (tested using very small datasets) | No | Yes |
| [63] | NBC, RF, and SVM | Byte Code, API and file system | Not robust against noisy inputs | Yes | No |
| [97] | k-NN classification | Byte Code | Not robust against noisy inputs | Yes | No |

| \multicolumn{6}{c}{Continuation of Table 1} | | | | | |
|---|---|---|---|---|---|
| Work | Classification method | Features | Limitations | Scalability (Yes/No) | Robust against noisy inputs (Yes/No) |
| [104] | HMM | opcode sequences | Not robust against severe obfuscations techniques | Yes | Yes |
| [105] | HMM | mnemonic opcode sequences | Not robust against severe obfuscations techniques | Yes | Yes |
| [106] | HMM | opcode sequences | Not robust against severe obfuscations techniques | Yes | Yes |
| [9] | HMM | opcode sequences | Not robust against severe obfuscation techniques | Yes | Yes |
| \multicolumn{6}{c}{End of Table} | | | | | |

Table 2: Comparison Summary (SA: Similarity Analyzes; FA: Families Analysis; VA: Varients Analysis.

| \multicolumn{5}{c}{Begin of Table} | | | | |
|---|---|---|---|---|
| Paper | Detection | SA | FA | VA |
| Schultz et al [64] | ✓ | | | |
| Kolter and Maloof [65] | ✓ | | | |
| Ahmed et al. [96] | ✓ | | | |
| Chau et al. [103] | ✓ | | | |
| Firdausi et al. [50] | ✓ | | | |
| Anderson et al. [21] | ✓ | | | |
| Anderson et al. [41] | ✓ | | | |
| Eskandari et al. [85] | ✓ | | | |

| Continuation of Table 2 | | | | |
|---|---|---|---|---|
| **Paper** | **Detection** | **SA** | **FA** | **VA** |
| Santos et al. [23] | ✓ | | | |
| Vadrevu et al. [73] | ✓ | | | |
| Bai et al. [74] | ✓ | | | |
| Kruczkowski and Szynkiewicz [57] | ✓ | | | |
| Tamersoy et al. [75] | ✓ | | | |
| Uppal et al. [58] | ✓ | | | |
| Chen et al. [78] | ✓ | | | |
| Ghiasi et al. [59] | ✓ | | | ✓ |
| Kwon et al. [98] | ✓ | | | |
| Mao et al. [99] | ✓ | | | |
| Saxe and Berlin [25] | ✓ | | | |
| Wuchner et al. [63] | ✓ | | | |
| Raff and Nicholas [97] | ✓ | | ✓ | |
| Gharacheh et al.[79] | | | | ✓ |
| Khodamoradi et al. [80] | | | | ✓ |
| Upchurch et al. [83] | | | | ✓ |
| Liang et al. [46] | | | | ✓ |
| Vadrevu and Perdisci [47] | | | | ✓ |
| Huang et al. [101] | | | ✓ | |
| Park et al. [51] | | | ✓ | |
| Ye et al. [102] | | | ✓ | |
| Dahl et al. [24] | | | ✓ | |
| Hu et al. [70] | | | ✓ | |
| Islam et al. [22] | | | ✓ | |
| Kong and Yan [71] | | | ✓ | |
| Nari and Ghorbani[55] | | | ✓ | |
| Ahmadi et al. [76] | | | ✓ | |
| Lin et al. [61] | | | ✓ | |
| Kawaguchi and Omote [60] | | | ✓ | |
| Mohaisen et al. [62] | | | ✓ | |
| Pai et al. [133] | | ✓ | | |
| Bailey et al. [48] | | ✓ | | |
| Bayer et al. [49] | | ✓ | | |
| Chen et al. [14] | | | ✓ | |

| Continuation of Table 2 | | | | |
|---|---|---|---|---|
| **Paper** | **Detection** | **SA** | **FA** | **VA** |
| Cesare and Xiang [15] | | | ✓ | |
| Anderson et al. [41] | | | ✓ | |
| Cordy et al. [93] | | | ✓ | |
| Fredrikson et al. [43] | | | ✓ | |
| Rieck et al. [53] | | ✓ | | |
| Palahan et al. [56] | | ✓ | | |
| Santos et al. [72] | | ✓ | | |
| Egele et al. [128] | | ✓ | | |
| Kolter and Maloof [17] | ✓ | | | |
| Moskovitch et al. [18] | ✓ | | | |
| End of Table | | | | |

## 6. Challenges and Issues

Based on the characterization explained in Section 5, we discuss here the challenges and/or issues of the surveyed articles.

### 6.1. Malware Evading Techniques

In this section, we introduce the common techniques that are used by malware authors to evade detection.

#### 6.1.1. Obfuscation

The term of obfuscation mainly refers to the techniques that are used to create a variant of the original code without affecting its functionality. The purpose of obfuscation is usually to hide the real logic of the original code or to evade signature-based detector or function clone detector. A few commonly used obfuscation techniques are as follows:

1. Dead-Code Insertion [13]: insert useless instructions (e.g., nop) or insert some instructions that only affect unused variables.
2. Code Transposition [13]: change the order of the independent instructions.
3. Register Reassignment [13]: exchange the usage of registers for the storage of data/address in a specific live range.
4. Instruction Substitution [13]: replace an instruction with equivalent instructions.

5. Control Flow Flattening [134]: 1) break up the body of the function to basic blocks 2) put all basic blocks which were originally at different nesting levels next to each other 3) encapsulate the basic blocks in a selective structure (a switch statement in the C++) 4) encapsulate the selection in a loop.

6. Bogus Control Flow [135]: for a basic block, add a new basic block which contains an opaque predicate and then make a conditional jump to the original basic block.

### 6.1.2. Packing

Packing is a technique to compress/encrypt an executable, where those packed files will be uncompressed/decrypted during runtime. It means that a static analyzer cannot see the real code since it doesn't run the executable. Packing is used not only for malware but also for the protection of Goodware schemes [15, 41]. According to the statistics conducted by Anderson et al. [41], 47.56% of the malware are packed and 19.59% of the Goodware are packed in their dataset.

### 6.1.3. Polymorphism

Polymorphism is also a technique that is based on encryption and decryption. A polymorphic malware contains two parts: the polymorphism engine and the real program which performs the malicious functions. The former mutates the encryption algorithms and keys when it replicates and the code of the latter per se is fixed but it is encrypted by the former in different ways during runtime. This way, the whole polymorphic malware program would look different at each generation [136].

### 6.1.4. Metamorphism

A metamorphic malware re-programs itself when it replicates. Consequently, in each generation, the whole program body is modified using code obfuscation techniques while the functionality is kept unchanged [136]. Metamorphic malware is considered to be more difficult to write than polymorphic malware.

### 6.2. Adversarial Attack and Defense

Since the direction of the recent research is to automate the process of malware analysis using machine learning techniques, the proposed solutions should be robust against adversarial examples, which are inputs designed

by an attacker to fool the machine learning models and make it generate erroneous decisions (e.g., making the malware analysis tools unable to detect malicious code). It has been recently shown that machine learning models, including deep neural networks, are quite vulnerable to adversarial examples. It is easy for an attacker to create "adversarial examples" [137] to fool a machine learning model through simply perpetuating parts of the inputs.

*6.2.1. Adversarial Attack*

Adversarial samples are crafted from normal samples with minimum perturbations on input variables to confuse a classifier without breaking the functionality of the original samples. It is natural that the perturbations should be based on the derivative of the loss function with respect to the classifier's input variables since derivatives show the directions of changes on the input that is the most effective for changing the output. So a differentiable classifier is required to create adversarial samples and deep learning models are just differentiable and effective classifiers. Studies show that adversarial samples generated to fool one model can fool a totally different model [138, 139]. Therefore, as deep learning models are proposed for the malware detection field, malware authors have better opportunities to craft adversarial examples to evade the detection of any machine learning models.

A formal description of the problem to craft an adversarial $x^*$ to be misclassified by a classifier $f$ is

$$min \ ||\delta_x|| \tag{5}$$
$$s.t. \ \ x^* = x + \delta_x, f(x^*) \neq f(x) \tag{6}$$

where $|| \cdot ||$ can be any norm and x is the sample to be perturbed.

Goodfellow et al. [140] present a fast gradient sign method in which the adversarial perturbation is determined by multiplying the gradients' sign of the sample $S$ with some coefficient to control the scale of perturbation. Papernot et al. [141] propose a forward derivative method which evaluates the sensitivity of the output to each input component using its Jacobian matrix and then constructs adversarial saliency maps based on the Jacobian matrix, indicating which input features to be included in the perturbation.

Compared with perturbing an adversarial image sample, there are some constraints on perturbing a malware sample since most of the features of malware are discrete rather than real-valued and the functionality should be intact. Thus, previous methods for perturbation of real-valued features need

to be adapted and some binary features can not be changed from "1" to "0" since "1" means that the feature exists and that the change in this direction may break the functionality.

Grosse et al. [28] propose a technique to craft adversarial Android malware. Inspired by Papernot et al. [141], [28] use the Jacobian matrix to examine which features have the greatest potential to lead to the prediction of a malicious program as being Goodware. They only allow distortions to no more than 20 features. All the features are binary features. To maintain the functionality of the adversarial example, they add two constraints: 1) only adjust manifest features that relate to the AndroidManifest.xml file. This file is available in any Android application; 2) it should be done by adding a single line of code to it. Using their method, a state-of-the-art feed-forward neural network which achieves 98% of accuracy on the original dataset is misled by 63% of the adversarial malware samples.

### 6.2.2. Adversarial Defense

Grosse et al. [28] try two methods to defend against adversarial attack. The first is to apply distillation [142, 141] to counter adversarial samples, which successfully reduces misclassification rate by 38.5% in some case. The second is adversarial training [140] which consists of training the model on the original dataset and then training the model again only on the adversarial samples for a few epochs. The misclassification rate is reduced to 67% from 73% through adversarial training.

Wang et al. [29] defend against adversarial attacks by randomly nullifying input features. Their nullification is similar to dropout since in both mechanisms some input features are randomly set to 0. The main difference with dropout is that the model don't drop any input feature during the test but in nullification some features are still dropped randomly during the test. Specifically, for each sample in any dataset, a nullification rate is sampled under a Gaussian distribution and the dimensions (features) to drop are sampled uniformly. The intuition is that nullification makes their architecture non-deterministic so that the attackers can't examine the importance of features and so it's hard for them to detect and exploit the "blind spots" of classifiers. In their experiments, the features are the invoked windows system DLL files and they use Jacobian-based saliency map to pick up to 10 features for each sample to perturb. Experimental results show that their method can improve the resistance to adversarial samples and that the best resistance is 64.86% and is achieved with a nullification rate of 10%. How-

ever, a theoretic problem of their approach is when adversarial samples are cross-model [138, 139]. Thus, even though nullification can harm the ability of an adversary to use this model to craft adversarial samples, the adversary can use other models (i.e., the same neural network without nullification) to craft adversarial samples which can also evade the one equipped with nullification. Therefore, there is no theoretic proof or evidence to show whether nullification can improve the resistance against adversarial samples crafted from other deep learning models.

## 6.3. Efficiency and Scalability

A practical malware search engine can help security engineers obtain malware search results on-the-fly when they are making analysis. Instant feedback provides the engineer the structure of a given malware that is under investigation [92]. One should note that scalability is an important factor as the number of malware in the database needs to scale up to millions. It is also a critical issue for producing a reliable malware search engine. For practical applications, a malware search engine' efficiency and scalability should be evaluated using a large repository in order to measure both its accuracy and latency.

## 7. Research Direction

The above contributions are effective in addressing some interesting research gaps in the literature. However, some points still need further study and investigation. The following research avenues could be further explored based on our literature review:

## 7.1. Robust Solutions

Although the discussed solutions in the literature review have paved the road for a reliable Malware Detection System (MDS) through extracting robust and useful features, the solution still needs to reduce human interaction. Thus, an automated system is required to take the data and automatically abstract and extract robust features from them. For this purpose, deep learning techniques could be the best candidate to replace the existing feature extraction approaches. The solution can be designed and implemented using different Deep Learning architectures (e.g., Generative Adversarial Networks, Stacked Denosing Autoencoder, Restricted Boltzman Machine, and

Variational Autoencoder) for auto-abstraction and extraction of robust features to significantly enhance the detection under heterogeneous, changing and noisy environments.

Recently, Ding et al. [143] propose a robust and accurate assembly clone search platform named Asm2Vec. The proposed platform enables engineers to automatically learns a vector representation of any assembly function by discriminating it from others functions. Also, the platform allows engineers to jointly learn the semantic relationships of assembly functions based on assembly code [143]. This, in turn enables us to construct useful and robust features to make efficient and reliable assembly clone search. The proposed learning representation is inspired by the Distributed Memory Model of Paragraph Vectors (PV-DM) model, which is used to learn a vectorized representation of a text paragraph [144]. The PV-DM model is fundamentally based on Word2Vec [145], which is used to learn vector representation of words. This is done by enabling words with similar meaning to be mapped to a similar position in the vector space. For example, "good" and "great" are close to each other, whereas "great" and "Japan" are more distant. Learning the vector representation of words becomes possible thanks to the concept of Distributed Vector Representation (DVR) of words, a well known method used for learning the word vectors. In particular, DVS exploits the power of machine learning models (usually Neural Networks) by training machine learning models to predict a word (i.e., target word) given the other words in a context. In the process of predicting the target word, we learn the vector representation of the target word.

The PV-DM model is inspired by Word2Vec by using the idea for learning the word vectors. In the PV-DM model, both word vectors and paragraph vectors are asked to contribute to the prediction of the target word given many contexts sampled from the paragraph [144]. This process (i.e., predicting the target word) allows us to learn the vector representation of the paragraph. Ding et al. [143] exploit the power of the PV-DM model to learn the vector representation of assembly functions based on assembly code. This is done by mapping assembly function (i.e., repository function) and the function's input tokens (i.e., instructions) to a unique vector. The machine learning model is then trained to predict a target token given the function and its tokens in a context. This process enables us to learn the vector representation of the function.

In fact, the solution should be able not only to accommodate unknown variants of known malware but also to accommodate unknown variants of

unknown malware. These solutions should also be robust against adversarial attacks. Although some works have already addressed this problem, these solutions are mostly based on adversarial training [146] and are not mature enough to combine the extraction of robust and useful features to protect the system against adversarial examples. Thus, the solution should not only be robust against complex and noisy data but also against adversarial examples.

## 7.2. Collaborative Solutions

Computer and communication systems are becoming more and more complex and vulnerable to intrusions. Cyber attacks are also becoming more complex and harder to analyse and recognize. In fact, it became increasingly difficult for a single MDS to recognize all intrusions, because of limited knowledge about the evolution of malware. The recent works in intrusion detection and malware analysis [147, 148, 149] have shown experimentally that the detection accuracy can be significantly improved, compared to the traditional single MDS, when MDSs cooperate with each other. In collaborative environment, each MDS can consult other MDSs about suspicious malware to increase the decision accuracy. Figure 4 shows an example of cooperative MDS.

Figure 4: The proposed taxonomy

Recently, Man and Huh [147] and Singh et al. [148] design a collaborative MDS, which enables malware-detection-alerts to be exchanged from different distributed detectors. Moreover, knowledge are enabled to be exchanged between nodes. In addition, Dermott et al. [150] propose a collaborative MDS in a cloud-computing environment. The proposed framework use the Dempster-Shafer theory of evidence [151] in order to combine the decisions form different malware detectors. The received decisions are aggregated to take the final decision regarding a suspicious malware. This technique has a shortcoming: its centralized-based architecture, whereby a reliable third-party is used for combining feedback and coordinating MDS.

In fact, the design of a cooperative MDS should take into consideration the following three properties (challenges): trustworthiness, fairness and sustainability. By trustworthiness, we mean that the MDS should be able to ensure that it will consult, cooperate and share knowledge with trusted parties (i.e., MDSs). By fairness, we mean that the MDS should be able to guarantee that mutual benefits will be achieved through minimizing the chance of cooperating with selfish MDSs. This is useful to give MDSs the motivation to

42

participate in the community. Finally, by sustainability, we mean enabling an MDS to proactively take decisions about suspicious attacks, regardless if the complete feedback have been received from consulted MDSs or not. Thus, the proposed solution will be applicable in real-time environments, where MDSs should take decisions about suspicious malware quickly.

## 7.3. Sustainable Solutions

The power of most malware analysis tools is largely based on the amount of knowledge that they have about Malware and dangerous attacks. In fact, supervised machine learning algorithms such as SVM, used by MDS, are heavily dependent on labeled data to learn how to effectively classify malicious and normal behaviours [152]. However, obtaining data on malicious behaviours is challenging and dangerous, especially if we are required to launch real attacks on production systems and put users, applications and systems at risk. To address this problem, we may need to have an efficient approach to synthesize new malware and augment our training data, in order to improve machine learning-based MDSs.

Generative models such as Generative adversarial Networks (GANs) [153] can be used to generate synthetic malware and enhance the detection accuracy of machine learning-based MDS, by augmenting Malware training sets. We encourage researchers to investigate the use of GANs, which have shown unprecedented ability in generating high quality new synthetic data, to generate malware variants. In particular, they need to design new algorithms to effectively and efficiently train GANs on the existing malware that are available in the repository in order to learn how to generate variants of them. To this end, researchers are required to collect a large volume of malware samples that consists of different attributes (vulnerabilities, targeted users, targeted hosts, etc.) from the public domain. Since GANs are only defined for real-valued, continued data and the design of malware is based on sequences of discrete tokens (bytes), special extensions should be applied on the original GANs theory. For example, we may need to integrate GANs with recurrent neural networks (RNNs) to tackle the problem of sequenced data [154]. Moreover, to address the problem of discrete data, we may need to place in parallel a dense layer per categorical variable, followed by Gumbel-Softmax activation and a concatenation to get the final output [155].

## 8. Conclusion

In this paper, we provide a comprehensive survey on publications that contributed to malware classification and composition analysis. There are four main contributions in our work. First, we proposed an organization of reviewed paper according to three dimensions: the purpose of the analysis (malware classification or composition analysis), the type of features obtained from samples, and the algorithms used to manipulate these features. Second, we provided a comparative analysis of the existing malware classification and composition analysis techniques, while structuring them according to the proposed taxonomy. Third, We determined the main issues and challenges associated with malware classification and composition analysis. Finally, we identified a number of emergent topics in the discussed field, such as collaborative malware analysis system, with guidelines on how to improve solutions to address the new challenges.

The above contributions are effective in addressing some interesting research gaps in the literature. However, some points still need further study and investigation. The following research avenues could be further explored in order to achieve better accuracy and efficient solutions compared to the state-of-the-art. The first avenue is the design of cooperative MDS to address the problem of limited and incomplete knowledge about malware. Through collaboration, an MDS can consult other MDSs about suspicious malware and increase the decision accuracy. To this end, we identify three challenges that should be addressed in cooperative MDS: trustworthiness, fairness and sustainability. Second, the design of robust MDS by enabling the automatic extraction of robust features from samples. The solution should be able not only to accommodate unknown variants of known malware but also to accommodate unknown variants of unknown malware. Moreover, the solution should be robust against adversarial attacks. Finally, the design of sustainable MDS by enabling an MDS to synthetically generate new malicious and benign code in order to enhance the accuracy of machine learning-based malware classification methods.

## Acknowledgments

## References

[1] Malware statistics and facts for 2020, 2020 (accessed 2020-03-17). https://www.comparitech.com/antivirus/malware-statistics-facts/.

[2] Malware Numbers 2017, 2019 (Accessed: 2019-08-17). https://www.gdatasoftware.com/blog/2018/03/30610-malware-number-2017.

[3] G. Suarez-Tangil, J. E. Tapiador, P. Peris-Lopez, A. Ribagorda, Evolution, detection and analysis of malware for smart devices, IEEE Communications Surveys & Tutorials 16 (2013) 961–987.

[4] J. P. Tailor, A. D. Patel, A comprehensive survey: ransomware attacks prevention, monitoring and damage control, Int. J. Res. Sci. Innov 4 (2017) 116–121.

[5] B. Vignau, R. Khoury, S. Hallé, 10 years of iot malware: A feature-based taxonomy, in: 2019 IEEE 19th International Conference on Software Quality, Reliability and Security Companion (QRS-C), IEEE, 2019, pp. 458–465.

[6] Z. Xu, H. Wang, Z. Xu, X. Wang, Power attack: An increasing threat to data centers., in: NDSS, 2014.

[7] K. Kimani, V. Oduol, K. Langat, Cyber security challenges for iot-based smart grid networks, International Journal of Critical Infrastructure Protection 25 (2019) 36–49.

[8] M. Jakobsson, Z. Ramzan, Crimeware: understanding new attacks and defenses, Addison-Wesley Professional, 2008.

[9] W. Wong, M. Stamp, Hunting for metamorphic engines, Journal in Computer Virology 2 (2006) 211–229.

[10] N. Tariq, Impact of cyberattacks on financial institutions, Journal of Internet Banking and Commerce 23 (2018) 1–11.

45

[11] L. Chen, Y. Ye, T. Bourlai, Adversarial machine learning in malware detection: Arms race between evasion attack and defense, in: 2017 European Intelligence and Security Informatics Conference (EISIC), IEEE, 2017, pp. 99–106.

[12] M. G. Schultz, E. Eskin, F. Zadok, S. J. Stolfo, Data mining methods for detection of new malicious executables, in: Security and Privacy, 2001. S&P 2001. Proceedings. 2001 IEEE Symposium on, IEEE, 2001, pp. 38–49.

[13] M. Christodorescu, S. Jha, Static analysis of executables to detect malicious patterns, Technical Report, WISCONSIN UNIV-MADISON DEPT OF COMPUTER SCIENCES, 2006.

[14] J. Chen, M. H. Alalfi, T. R. Dean, Y. Zou, Detecting android malware using clone detection, Journal of Computer Science and Technology 30 (2015) 942–956.

[15] S. Cesare, Y. Xiang, Classification of malware using structured control flow, in: Proceedings of the Eighth Australasian Symposium on Parallel and Distributed Computing-Volume 107, Australian Computer Society, Inc., 2010, pp. 61–70.

[16] Y. Ye, T. Li, K. Huang, Q. Jiang, Y. Chen, Hierarchical associative classifier (hac) for malware detection from the large and imbalanced gray list, Journal of Intelligent Information Systems 35 (2010) 1–20.

[17] J. Z. Kolter, M. A. Maloof, Learning to detect malicious executables in the wild, in: Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining, ACM, 2004, pp. 470–478.

[18] R. Moskovitch, C. Feher, N. Tzachar, E. Berger, M. Gitelman, S. Dolev, Y. Elovici, Unknown malcode detection using opcode representation, in: Intelligence and Security Informatics, Springer, 2008, pp. 204–215.

[19] J. Dai, R. K. Guha, J. Lee, Efficient virus detection using dynamic instruction sequences., JCP 4 (2009) 405–414.

[20] L. Nataraj, S. Karthikeyan, G. Jacob, B. Manjunath, Malware images: visualization and automatic classification, in: Proceedings of the

8th international symposium on visualization for cyber security, ACM, 2011, p. 4.

[21] B. Anderson, D. Quist, J. Neil, C. Storlie, T. Lane, Graph-based malware detection using dynamic analysis, Journal in computer Virology 7 (2011) 247–258.

[22] R. Islam, R. Tian, L. M. Batten, S. Versteeg, Classification of malware based on integrated static and dynamic features, Journal of Network and Computer Applications 36 (2013) 646–656.

[23] I. Santos, J. Devesa, F. Brezo, J. Nieves, P. G. Bringas, Opem: A static-dynamic approach for machine-learning-based malware detection, in: International Joint Conference CISIS'12-ICEUTE 12-SOCO 12 Special Sessions, Springer, 2013, pp. 271–280.

[24] G. E. Dahl, J. W. Stokes, L. Deng, D. Yu, Large-scale malware classification using random projections and neural networks, in: Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on, IEEE, 2013, pp. 3422–3426.

[25] J. Saxe, K. Berlin, Deep neural network based malware detection using two dimensional binary program features, in: Malicious and Unwanted Software (MALWARE), 2015 10th International Conference on, IEEE, 2015, pp. 11–20.

[26] W. Huang, J. W. Stokes, Mtnet: a multi-task neural network for dynamic malware classification, in: International Conference on Detection of Intrusions and Malware, and Vulnerability Assessment, Springer, 2016, pp. 399–418.

[27] B. Kolosnjaji, A. Zarras, G. Webster, C. Eckert, Deep learning for classification of malware system call sequences, in: Australasian Joint Conference on Artificial Intelligence, Springer, 2016, pp. 137–149.

[28] K. Grosse, N. Papernot, P. Manoharan, M. Backes, P. McDaniel, Adversarial examples for malware detection, in: European Symposium on Research in Computer Security, Springer, 2017, pp. 62–79.

[29] Q. Wang, W. Guo, K. Zhang, A. G. Ororbia II, X. Xing, X. Liu, C. L. Giles, Adversary resistant deep neural networks with an application

1187 to malware detection, in: Proceedings of the 23rd ACM SIGKDD
1188 International Conference on Knowledge Discovery and Data Mining,
1189 ACM, 2017, pp. 1145–1153.

[30] D. Ucci, L. Aniello, R. Baldoni, Survey of machine learning techniques
for malware analysis, Computers & Security (2018).

[31] M. K. Sahu, M. Ahirwar, A. Hemlata, A review of malware detection
based on pattern matching technique, Int. J. of Computer Science and
Information Technologies (IJCSIT) 5 (2014) 944–947.

[32] A. Souri, R. Hosseini, A state-of-the-art survey of malware detection
approaches using data mining techniques, Human-centric Computing
and Information Sciences 8 (2018) 3.

[33] Z. Bazrafshan, H. Hashemi, S. M. H. Fard, A. Hamzeh, A survey on
heuristic malware detection techniques, in: The 5th Conference on
Information and Knowledge Technology, IEEE, 2013, pp. 113–120.

[34] A. Shabtai, R. Moskovitch, Y. Elovici, C. Glezer, Detection of mali-
cious code by applying machine learning classifiers on static features: A
state-of-the-art survey, information security technical report 14 (2009)
16–29.

[35] I. Basu, N. Sinha, D. Bhagat, S. Goswami, Malware detection based
on source data using data mining: A survey, American Journal Of
Advanced Computing 3 (2016) 18–37.

[36] Y. Ye, T. Li, D. Adjeroh, S. S. Iyengar, A survey on malware detection
using data mining techniques, ACM Computing Surveys (CSUR) 50
(2017) 41.

[37] O. Or-Meir, N. Nissim, Y. Elovici, L. Rokach, Dynamic malware anal-
ysis in the modern era—a state of the art survey, ACM Computing
Surveys (CSUR) 52 (2019) 88.

[38] J. Barriga, S. Yoo, Malware detection and evasion with machine learn-
ing techniques: A survey, International Journal of Applied Engineering
Research 12 (2017).

[39] A. Damodaran, F. Di Troia, C. A. Visaggio, T. H. Austin, M. Stamp, A comparison of static, dynamic, and hybrid analysis for malware detection, Journal of Computer Virology and Hacking Techniques 13 (2017) 1–12.

[40] U. Bayer, A. Moser, C. Kruegel, E. Kirda, Dynamic analysis of malicious code, Journal in Computer Virology 2 (2006) 67–77.

[41] B. Anderson, C. Storlie, T. Lane, Improving malware classification: bridging the static/dynamic gap, in: Proceedings of the 5th ACM workshop on Security and artificial intelligence, ACM, 2012, pp. 3–14.

[42] P. Royal, M. Halpin, D. Dagon, R. Edmonds, W. Lee, Polyunpack: Automating the hidden-code extraction of unpack-executing malware, in: Computer Security Applications Conference, 2006. ACSAC'06. 22nd Annual, IEEE, 2006, pp. 289–300.

[43] M. Fredrikson, S. Jha, M. Christodorescu, R. Sailer, X. Yan, Synthesizing near-optimal malware specifications from suspicious behaviors, in: Security and Privacy (SP), 2010 IEEE Symposium on, IEEE, 2010, pp. 45–60.

[44] U. A. Force, Analysis of the intel pentium's ability to support a secure virtual machine monitor, in: Proceedings of the 9th USENIX Security Symposium., 2000, p. 129.

[45] J. Rutkowska, Redpill: Detect vmm using (almost) one cpu instruction, http://invisiblethings. org/papers/redpill. html (2004).

[46] G. Liang, J. Pang, C. Dai, A behavior-based malware variant classification technique, International Journal of Information and Education Technology 6 (2016) 291.

[47] P. Vadrevu, R. Perdisci, Maxs: Scaling malware execution with sequential multi-hypothesis testing, in: Proceedings of the 11th ACM on Asia Conference on Computer and Communications Security, ACM, 2016, pp. 771–782.

[48] M. Bailey, J. Oberheide, J. Andersen, Z. M. Mao, F. Jahanian, J. Nazario, Automated classification and analysis of internet malware,

1248      in: International Workshop on Recent Advances in Intrusion Detection,
1249      Springer, 2007, pp. 178–197.

1250 [49] U. Bayer, P. M. Comparetti, C. Hlauschek, C. Kruegel, E. Kirda, Scal-
1251      able, behavior-based malware clustering., in: NDSS, volume 9, Cite-
1252      seer, 2009, pp. 8–11.

1253 [50] I. Firdausi, A. Erwin, A. S. Nugroho, et al., Analysis of machine
1254      learning techniques used in behavior-based malware detection, in: 2010
1255      second international conference on advances in computing, control, and
1256      telecommunication technologies, IEEE, 2010, pp. 201–203.

1257 [51] Y. Park, D. Reeves, V. Mulukutla, B. Sundaravel, Fast malware clas-
1258      sification by automated behavioral graph matching, in: Proceedings of
1259      the Sixth Annual Workshop on Cyber Security and Information Intel-
1260      ligence Research, ACM, 2010, p. 45.

1261 [52] M. Lindorfer, C. Kolbitsch, P. M. Comparetti, Detecting environment-
1262      sensitive malware, in: International Workshop on Recent Advances in
1263      Intrusion Detection, Springer, 2011, pp. 338–357.

1264 [53] K. Rieck, P. Trinius, C. Willems, T. Holz, Automatic analysis of mal-
1265      ware behavior using machine learning, Journal of Computer Security
1266      19 (2011) 639–668.

1267 [54] P. M. Comar, L. Liu, S. Saha, P.-N. Tan, A. Nucci, Combining super-
1268      vised and unsupervised learning for zero-day malware detection, in:
1269      2013 Proceedings IEEE INFOCOM, IEEE, 2013, pp. 2022–2030.

1270 [55] S. Nari, A. A. Ghorbani, Automated malware classification based on
1271      network behavior, in: 2013 International Conference on Computing,
1272      Networking and Communications (ICNC), IEEE, 2013, pp. 642–647.

1273 [56] S. Palahan, D. Babić, S. Chaudhuri, D. Kifer, Extraction of statistically
1274      significant malware behaviors, in: Proceedings of the 29th Annual
1275      Computer Security Applications Conference, ACM, 2013, pp. 69–78.

1276 [57] M. Kruczkowski, E. N. Szynkiewicz, Support vector machine for
1277      malware analysis and classification, in: Proceedings of the 2014
1278      IEEE/WIC/ACM International Joint Conferences on Web Intelligence

1279  (WI) and Intelligent Agent Technologies (IAT)-Volume 02, IEEE Com-
1280  puter Society, 2014, pp. 415–420.

1281  [58] D. Uppal, R. Sinha, V. Mehra, V. Jain, Malware detection and clas-
1282  sification based on extraction of api sequences, in: 2014 International
1283  Conference on Advances in Computing, Communications and Infor-
1284  matics (ICACCI), IEEE, 2014, pp. 2337–2342.

1285  [59] M. Ghiasi, A. Sami, Z. Salehi, Dynamic vsa: a framework for mal-
1286  ware detection based on register contents, Engineering Applications of
1287  Artificial Intelligence 44 (2015) 111–122.

1288  [60] N. Kawaguchi, K. Omote, Malware function classification using apis in
1289  initial behavior, in: 2015 10th Asia Joint Conference on Information
1290  Security, IEEE, 2015, pp. 138–144.

1291  [61] C.-T. Lin, N.-J. Wang, H. Xiao, C. Eckert, Feature selection and
1292  extraction for malware classification., J. Inf. Sci. Eng. 31 (2015) 965–
1293  992.

1294  [62] A. Mohaisen, O. Alrawi, M. Mohaisen, Amal: High-fidelity, behavior-
1295  based automated malware analysis and classification, computers &
1296  security 52 (2015) 251–266.

1297  [63] T. Wüchner, M. Ochoa, A. Pretschner, Robust and effective malware
1298  detection through quantitative data flow graph metrics, in: Interna-
1299  tional Conference on Detection of Intrusions and Malware, and Vulner-
1300  ability Assessment, Springer, 2015, pp. 98–118.

1301  [64] M. G. Schultz, E. Eskin, F. Zadok, S. J. Stolfo, Data mining methods
1302  for detection of new malicious executables, in: Proceedings 2001 IEEE
1303  Symposium on Security and Privacy. S&P 2001, IEEE, 2000, pp. 38–49.

1304  [65] J. Z. Kolter, M. A. Maloof, Learning to detect and classify malicious
1305  executables in the wild, Journal of Machine Learning Research 7 (2006)
1306  2721–2744.

1307  [66] S. Attaluri, S. McGhee, M. Stamp, Profile hidden markov models and
1308  metamorphic virus detection, Journal in computer virology 5 (2009)
1309  151–169.

[67] M. Siddiqui, M. C. Wang, J. Lee, Detecting internet worms using data mining techniques, Journal of Systemics, Cybernetics and Informatics 6 (2009) 48–53.

[68] I. Santos, J. Nieves, P. G. Bringas, Semi-supervised learning for unknown malware detection, in: International Symposium on Distributed Computing and Artificial Intelligence, Springer, 2011, pp. 415–422.

[69] Z. Chen, M. Roussopoulos, Z. Liang, Y. Zhang, Z. Chen, A. Delis, Malware characteristics and threats on the internet ecosystem, Journal of Systems and Software 85 (2012) 1650–1672.

[70] X. Hu, K. G. Shin, S. Bhatkar, K. Griffin, Mutantx-s: Scalable malware clustering based on static features, in: Proceedings of the USENIX Annual Technical Conference (ATC), 2013, pp. 187–198.

[71] D. Kong, G. Yan, Discriminant malware distance learning on structural information for automated malware classification, in: Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining, ACM, 2013, pp. 1357–1365.

[72] I. Santos, F. Brezo, X. Ugarte-Pedrero, P. G. Bringas, Opcode sequences as representation of executables for data-mining-based unknown malware detection, Information Sciences 231 (2013) 64–82.

[73] P. Vadrevu, B. Rahbarinia, R. Perdisci, K. Li, M. Antonakakis, Measuring and detecting malware downloads in live network traffic, in: European Symposium on Research in Computer Security, Springer, 2013, pp. 556–573.

[74] J. Bai, J. Wang, G. Zou, A malware detection scheme based on mining format information, The Scientific World Journal 2014 (2014).

[75] A. Tamersoy, K. Roundy, D. H. Chau, Guilt by association: large scale malware detection by mining file-relation graphs, in: Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining, ACM, 2014, pp. 1524–1533.

[76] M. Ahmadi, D. Ulyanov, S. Semenov, M. Trofimov, G. Giacinto, Novel feature extraction, selection and fusion for effective malware family

classification, in: Proceedings of the sixth ACM conference on data and application security and privacy, ACM, 2016, pp. 183–194.

[77] A. Caliskan-Islam, R. Harang, A. Liu, A. Narayanan, C. Voss, F. Yamaguchi, R. Greenstadt, De-anonymizing programmers via code stylometry, in: Proceedings of the 24th USENIX Security Symposium, 2015, pp. 255–270.

[78] L. Chen, T. Li, M. Abdulhayoglu, Y. Ye, Intelligent malware detection based on file relation graphs, in: Proceedings of the 2015 IEEE 9th International Conference on Semantic Computing (IEEE ICSC 2015), IEEE, 2015, pp. 85–92.

[79] M. Gharacheh, V. Derhami, S. Hashemi, S. M. H. Fard, Proposing an hmm-based approach to detect metamorphic malware, in: 2015 4th Iranian Joint Congress on Fuzzy and Intelligent Systems (CFIS), IEEE, 2015, pp. 1–5.

[80] P. Khodamoradi, M. Fazlali, F. Mardukhi, M. Nosrati, Heuristic metamorphic malware detection based on statistics of assembly instructions using classification algorithms, in: 2015 18th CSI International Symposium on Computer Architecture and Digital Systems (CADS), IEEE, 2015, pp. 1–6.

[81] J. Sexton, C. Storlie, B. Anderson, Subroutine based detection of apt malware, Journal of Computer Virology and Hacking Techniques 12 (2016) 225–233.

[82] S. S. W. Piyanuntcharatsr, S. Adulkasem, C. Chantrapornchai, On the comparison of malware detection methods using data mining with two feature sets, International Journal of Security and Its Applications 9 (2015) 293–318.

[83] J. Upchurch, X. Zhou, Variant: a malware similarity testing framework, in: 2015 10th International Conference on Malicious and Unwanted Software (MALWARE), IEEE, 2015, pp. 31–39.

[84] J. Jang, D. Brumley, S. Venkataraman, Bitshred: feature hashing malware for scalable triage and semantic analysis, in: Proceedings of the 18th ACM conference on Computer and communications security, ACM, 2011, pp. 309–320.

[85] M. Eskandari, Z. Khorshidpour, S. Hashemi, Hdm-analyser: a hybrid analysis approach based on data mining techniques for malware detection, Journal of Computer Virology and Hacking Techniques 9 (2013) 77–93.

[86] M. Graziano, D. Canali, L. Bilge, A. Lanzi, E. Shi, D. Balzarotti, M. van Dijk, M. Bailey, S. Devadas, M. Liu, et al., Needles in a haystack: Mining information from public dynamic analysis sandboxes for malware intelligence, in: Proceedings of the 24th USENIX Security Symposium, 2015, pp. 1057–1072.

[87] A. Oliva, A. Torralba, Modeling the shape of the scene: A holistic representation of the spatial envelope, International journal of computer vision 42 (2001) 145–175.

[88] N. Bhodia, P. Prajapati, F. Di Troia, M. Stamp, Transfer learning for image-based malware classification, arXiv preprint arXiv:1903.11551 (2019).

[89] R. Agrawal, R. Srikant, et al., Fast algorithms for mining association rules, in: Proc. 20th int. conf. very large data bases, VLDB, volume 1215, 1994, pp. 487–499.

[90] C. Kruegel, E. Kirda, D. Mutz, W. Robertson, G. Vigna, Polymorphic worm detection using structural information of executables, in: International Workshop on Recent Advances in Intrusion Detection, Springer, 2005, pp. 207–226.

[91] S. H. H. Ding, B. C. M. Fung, P. Charland, Asm2vec: Boosting static representation robustness for binary clone search against code obfuscation and compiler optimization, in: 2019 IEEE Symposium on Security and Privacy (SP), IEEE, 2019, pp. 472–489.

[92] S. H. H. Ding, B. C. M. Fung, P. Charland, Kam1n0: Mapreduce-based assembly clone search for reverse engineering, in: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, ACM, 2016, pp. 461–470.

[93] J. R. Cordy, C. K. Roy, The nicad clone detector, in: Program Comprehension (ICPC), 2011 IEEE 19th International Conference on, IEEE, 2011, pp. 219–220.

[94] D. Baysa, R. M. Low, M. Stamp, Structural entropy and metamorphic malware, Journal of computer virology and hacking techniques 9 (2013) 179–192.

[95] R. Tian, L. M. Batten, S. Versteeg, Function length as a tool for malware classification, in: 2008 3rd International Conference on Malicious and Unwanted Software (MALWARE), IEEE, 2008, pp. 69–76.

[96] F. Ahmed, H. Hameed, M. Z. Shafiq, M. Farooq, Using spatio-temporal information in api calls with machine learning algorithms for malware detection, in: Proceedings of the 2nd ACM workshop on Security and artificial intelligence, ACM, 2009, pp. 55–62.

[97] E. Raff, C. Nicholas, An alternative to ncd for large sequences, lempel-ziv jaccard distance, in: Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, ACM, 2017, pp. 1007–1015.

[98] B. J. Kwon, J. Mondal, J. Jang, L. Bilge, T. Dumitraş, The dropper effect: Insights into malware distribution with downloader graph analytics, in: Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security, ACM, 2015, pp. 1118–1129.

[99] W. Mao, Z. Cai, D. Towsley, X. Guan, Probabilistic inference on integrity for access behavior based malware detection, in: International Symposium on Recent Advances in Intrusion Detection, Springer, 2015, pp. 155–176.

[100] M. Polino, A. Scorti, F. Maggi, S. Zanero, Jackdaw: Towards automatic reverse engineering of large datasets of binaries, in: International Conference on Detection of Intrusions and Malware, and Vulnerability Assessment, Springer, 2015, pp. 121–143.

[101] K. Huang, Y. Ye, Q. Jiang, Ismcs: an intelligent instruction sequence based malware categorization system, in: 2009 3rd International Conference on Anti-counterfeiting, Security, and Identification in Communication, IEEE, 2009, pp. 509–512.

[102] Y. Ye, T. Li, Y. Chen, Q. Jiang, Automatic malware categorization using cluster ensemble, in: Proceedings of the 16th ACM SIGKDD in-

ternational conference on Knowledge discovery and data mining, ACM, 2010, pp. 95–104.

[103] C. Nachenberg, J. Wilhelm, A. Wright, C. Faloutsos, Polonium: Tera-scale graph mining for malware detection (2010).

[104] A. Kalbhor, T. H. Austin, E. Filiol, S. Josse, M. Stamp, Dueling hidden markov models for virus analysis, Journal of Computer Virology and Hacking Techniques 11 (2015) 103–118.

[105] A. Raghavan, F. Di Troia, M. Stamp, Hidden markov models with random restarts versus boosting for malware detection, Journal of Computer Virology and Hacking Techniques 15 (2019) 97–107.

[106] C. Annachhatre, T. H. Austin, M. Stamp, Hidden markov models for malware classification, Journal of Computer Virology and Hacking Techniques 11 (2015) 59–73.

[107] S. J. Russell, P. Norvig, Artificial intelligence: a modern approach, Malaysia; Pearson Education Limited,, 2016.

[108] J. R. Quinlan, Induction of decision trees, Machine learning 1 (1986) 81–106.

[109] N. S. Altman, An introduction to kernel and nearest-neighbor non-parametric regression, The American Statistician 46 (1992) 175–185.

[110] B. E. Boser, I. M. Guyon, V. N. Vapnik, A training algorithm for optimal margin classifiers, in: Proceedings of the fifth annual workshop on Computational learning theory, ACM, 1992, pp. 144–152.

[111] F. V. Jensen, An introduction to Bayesian networks, volume 210, UCL press London, 1996.

[112] B. Liu, Y. Ma, C. K. Wong, Improving an association rule based classifier, in: European Conference on Principles of Data Mining and Knowledge Discovery, Springer, 2000, pp. 504–509.

[113] W. W. Cohen, Learning trees and rules with set-valued features, in: AAAI/IAAI, Vol. 1, 1996, pp. 709–716.

[114] L. K. Hansen, P. Salamon, Neural network ensembles, IEEE Transactions on Pattern Analysis & Machine Intelligence (1990) 993–1001.

[115] M. Pal, Random forest classifier for remote sensing classification, International Journal of Remote Sensing 26 (2005) 217–222.

[116] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, P.-A. Manzagol, Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion, Journal of machine learning research 11 (2010) 3371–3408.

[117] A. Ng, et al., Sparse autoencoder, CS294A Lecture notes 72 (2011) 1–19.

[118] O. Fink, E. Zio, U. Weidmann, Fuzzy classification with restricted boltzman machines and echo-state networks for predicting potential railway door system failures, IEEE Transactions on Reliability 64 (2015) 861–868.

[119] X. Yan, H. Cheng, J. Han, P. S. Yu, Mining significant graph patterns by leap search, in: Proceedings of the 2008 ACM SIGMOD international conference on Management of data, ACM, 2008, pp. 433–444.

[120] R. Wille, Restructuring lattice theory: an approach based on hierarchies of concepts, in: Ordered sets, Springer, 1982, pp. 445–470.

[121] P. Brémaud, Markov chains: Gibbs fields, Monte Carlo simulation, and queues, volume 31, Springer Science & Business Media, 2013.

[122] C. Kruegel, W. Robertson, F. Valeur, G. Vigna, Static disassembly of obfuscated binaries, in: USENIX Security Symposium, volume 13, 2004, pp. 18–18.

[123] C. Cifuentes, K. J. Gough, Decompilation of binary programs, Software: Practice and Experience 25 (1995) 811–829.

[124] C. Cifuentes, M. Van Emmerik, Uqbt: Adaptable binary translation at low cost, Computer 33 (2000) 60–66.

[125] H.-C. Shin, M. R. Orton, D. J. Collins, S. J. Doran, M. O. Leach, Stacked autoencoders for unsupervised feature learning and multiple

organ detection in a pilot study using 4d patient data, IEEE transactions on pattern analysis and machine intelligence 35 (2012) 1930–1943.

[126] Y.-l. Boureau, Y. L. Cun, et al., Sparse feature learning for deep belief networks, in: Advances in neural information processing systems, 2008, pp. 1185–1192.

[127] S. R. Eddy, Hidden markov models, Current opinion in structural biology 6 (1996) 361–365.

[128] M. Egele, M. Woo, P. Chapman, D. Brumley, Blanket execution: Dynamic similarity testing for program binaries and components, in: Proceedings of the 23rd USENIX Security Symposium, 2014, pp. 303–317.

[129] B. N. Narayanan, O. Djaneye-Boundjou, T. M. Kebede, Performance analysis of machine learning and pattern recognition algorithms for malware classification, in: 2016 IEEE National Aerospace and Electronics Conference (NAECON) and Ohio Innovation Summit (OIS), IEEE, 2016, pp. 338–342.

[130] T. M. Kebede, O. Djaneye-Boundjou, B. N. Narayanan, A. Ralescu, D. Kapp, Classification of malware programs using autoencoders based deep learning architecture and its application to the microsoft malware classification challenge (big 2015) dataset, in: 2017 IEEE National Aerospace and Electronics Conference (NAECON), IEEE, 2017, pp. 70–75.

[131] T. Messay-Kebede, B. N. Narayanan, O. Djaneye-Boundjou, Combination of traditional and deep learning based architectures to overcome class imbalance and its application to malware classification, in: NAECON 2018-IEEE National Aerospace and Electronics Conference, IEEE, 2018, pp. 73–77.

[132] V. S. P. Davuluru, B. N. Narayanan, E. J. Balster, Convolutional neural networks as classification tools and feature extractors for distinguishing malware programs, in: 2019 IEEE National Aerospace and Electronics Conference (NAECON), IEEE, 2019, pp. 273–278.

[133] S. Pai, F. Di Troia, C. A. Visaggio, T. H. Austin, M. Stamp, Clustering for malware classification, Journal of Computer Virology and Hacking Techniques 13 (2017) 95–107.

[134] T. László, Á. Kiss, Obfuscating c++ programs via control flow flattening, Annales Universitatis Scientarum Budapestinensis de Rolando Eötvös Nominatae, Sectio Computatorica 30 (2009) 3–19.

[135] Bogus Control Flow, 2020 (accessed 2020-03-10). https://github.com/obfuscator-llvm/obfuscator/wiki/Bogus-Control-Flow.

[136] X. Li, P. K. Loh, F. Tan, Mechanisms of polymorphic and metamorphic viruses, in: Intelligence and Security Informatics Conference (EISIC), 2011 European, IEEE, 2011, pp. 149–154.

[137] A. Kurakin, I. Goodfellow, S. Bengio, Adversarial examples in the physical world, arXiv preprint arXiv:1607.02533 (2016).

[138] J. Bruna, C. Szegedy, I. Sutskever, I. Goodfellow, W. Zaremba, R. Fergus, D. Erhan, Intriguing properties of neural networks (2013).

[139] N. Papernot, P. McDaniel, I. Goodfellow, Transferability in machine learning: from phenomena to black-box attacks using adversarial samples, arXiv preprint arXiv:1605.07277 (2016).

[140] I. J. Goodfellow, J. Shlens, C. Szegedy, Explaining and harnessing adversarial examples, CoRR abs/1412.6572 (2014).

[141] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, A. Swami, The limitations of deep learning in adversarial settings, in: Security and Privacy (EuroS&P), 2016 IEEE European Symposium on, IEEE, 2016, pp. 372–387.

[142] G. Hinton, O. Vinyals, J. Dean, Distilling the knowledge in a neural network, arXiv preprint arXiv:1503.02531 (2015).

[143] S. H. H. Ding, B. C. M. Fung, P. Charland, Asm2vec: Boosting static representation robustness for binary clone search against code obfuscation and compiler optimization, in: Proc. of the 40th International Symposium on Security and Privacy (S&P), IEEE Computer Society, San Francisco, CA, 2019, pp. 38–55.

[144] Q. Le, T. Mikolov, Distributed representations of sentences and documents, in: International conference on machine learning, 2014, pp. 1188–1196.

[145] T. Mikolov, K. Chen, G. Corrado, J. Dean, Efficient estimation of word representations in vector space, arXiv preprint arXiv:1301.3781 (2013).

[146] N. Carlini, D. Wagner, Audio adversarial examples: Targeted attacks on speech-to-text, in: 2018 IEEE Security and Privacy Workshops (SPW), IEEE, 2018, pp. 1–7.

[147] N. D. Man, E.-N. Huh, A collaborative intrusion detection system framework for cloud computing, in: Proceedings of the International Conference on IT Convergence and Security 2011, Springer, 2012, pp. 91–109.

[148] D. Singh, D. Patel, B. Borisaniya, C. Modi, Collaborative ids framework for cloud, International Journal of Network Security 18 (2016) 699–709.

[149] C. J. Fung, Q. Zhu, Facid: A trust-based collaborative decision framework for intrusion detection networks, Ad Hoc Networks 53 (2016) 17–31.

[150] A. Mac Dermott, Q. Shi, K. Kifayat, Collaborative intrusion detection in federated cloud environments, Journal of Computer Sciences and Applications 3 (2015) 10–20.

[151] G. Shafer, Dempster-shafer theory, Encyclopedia of artificial intelligence 1 (1992) 330–331.

[152] F. Pendlebury, F. Pierazzi, R. Jordaney, J. Kinder, L. Cavallaro, {TESSERACT}: Eliminating experimental bias in malware classification across space and time, in: Proceedings of the 28th USENIX Security Symposium), 2019, pp. 729–746.

[153] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative adversarial nets, in: Advances in neural information processing systems, 2014, pp. 2672–2680.

[154] D. J. Im, C. D. Kim, H. Jiang, R. Memisevic, Generating images with recurrent adversarial networks, arXiv preprint arXiv:1602.05110 (2016).

[155] E. Jang, S. Gu, B. Poole, Categorical reparameterization with gumbel-softmax, arXiv preprint arXiv:1611.01144 (2016).