

 Open access • Book Chapter • DOI:10.1007/978-3-642-04342-0\_7

## PE-Miner: Mining Structural Information to Detect Malicious Executables in Realtime

— [Source link](#) 

M. Zubair Shafiq, S. Momina Tabish, Fauzan Mirza, Muddassar Farooq

**Institutions:** National University of Computer and Emerging Sciences, University of the Sciences

**Published on:** 01 Oct 2009 - Recent Advances in Intrusion Detection

**Topics:** Malware

Related papers:

- [Data mining methods for detection of new malicious executables](#)
- [Learning to Detect and Classify Malicious Executables in the Wild](#)
- [McBoost: Boosting Scalability in Malware Collection and Analysis Using Statistical Classification of Executables](#)
- [PE-Probe: Leveraging Packer Detection and Structural Information to Detect Malicious Portable Executables](#)
- [Learning to detect malicious executables in the wild](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/pe-miner-mining-structural-information-to-detect-malicious-38jynfifmx>

# PE-Miner: Mining Structural Information to Detect Malicious Executables in Realtime

M. Zubair Shafiq<sup>1</sup>, S. Momina Tabish<sup>1,2</sup>, Fauzan Mirza<sup>2,1</sup>, Muddassar Farooq<sup>1</sup>

<sup>1</sup>Next Generation Intelligent Networks Research Center (nexGIN RC)  
National University of Computer & Emerging Sciences (FAST-NUCES)  
Islamabad, 44000, Pakistan

{zubair.shafiq,momina.tabish,muddassar.farooq}@nexginrc.org

<sup>2</sup>School of Electrical Engineering & Computer Science (SEECS)

National University of Sciences & Technology (NUST)

Islamabad, 44000, Pakistan

fauzan.mirza@seecs.edu.pk

**Abstract.** In this paper, we present an accurate and realtime PE-Miner framework that automatically extracts distinguishing features from portable executables (PE) to detect zero-day (i.e. previously unknown) malware. The distinguishing features are extracted using the structural information standardized by the Microsoft Windows operating system for executables, DLLs and object files. We follow a threefold research methodology: (1) identify a set of structural features for PE files which is computable in realtime, (2) use an efficient preprocessor for removing redundancy in the features' set, and (3) select an efficient data mining algorithm for final classification between benign and malicious executables.

We have evaluated PE-Miner on two malware collections, VX Heavens and Malfease datasets which contain about 11 and 5 thousand malicious PE files respectively. The results of our experiments show that PE-Miner achieves more than 99% detection rate with less than 0.5% false alarm rate for distinguishing between benign and malicious executables. PE-Miner has low processing overheads and takes only 0.244 seconds on the average to scan a given PE file. Finally, we evaluate the robustness and reliability of PE-Miner under several regression tests. Our results show that the extracted features are robust to different packing techniques and PE-Miner is also resilient to majority of crafty evasion strategies.

**Key words:** Data Mining, Malicious Executable Detection, Malware Detection, Portable Executables, Structural Information

## 1 Introduction

A number of non-signature based malware detection techniques have been proposed recently. These techniques mostly use heuristic analysis, behavior analysis, or a combination of both to detect malware. Such techniques are being actively investigated because of their ability to detect zero-day malware without any a

priori knowledge about them. Some of them have been integrated into the existing Commercial Off the Shelf Anti Virus (COTS AV) products, but have achieved only limited success [26], [13]. The most important shortcoming of these techniques is that they are not *realtime deployable*<sup>1</sup>. We, therefore, believe that the domain of *realtime deployable* non-signature based malware detection techniques is still open to novel research.

Non-signature based malware detection techniques are primarily criticized because of two inherent problems: (1) high *fp* rate, and (2) large processing overheads. Consequently, COTS AV products mostly utilize signature based detection schemes that provide low *fp* rate and have acceptable processing overheads. But it is a well-known fact that signature based malware detection schemes are unable to detect *zero-day* malware. We cite two reports to highlight the alarming rate at which new malware is proliferating. The first report is by Symantec that shows an increase of 468% in the number of malware from 2006 to 2007 [25]. The second report shows that the number of malware produced in 2007 alone was more than the total number of malware produced in the last 20 years [6]. These surveys suggest that signature based techniques cannot keep abreast with the security challenges of the new millennium because not only the size of the signatures' database will exponentially increase but also the time of matching signatures. These bottlenecks are even more relevant on resource constrained smart phones and mobile devices [3]. We, therefore, envision that in near future signature based malware detection schemes will not be able to meet the criterion of *realtime deployable* as well.

We argue that a malware detection scheme which is *realtime deployable* should use an intelligent yet simple static analysis technique. In this paper we propose a framework, called *PE-Miner*, which uses novel *structural features* to efficiently detect malicious PE files. PE is a file format which is standardized by the Microsoft Windows operating systems for executables, dynamically linked libraries (DLL), and object files. We follow a threefold research methodology in our static analysis: (1) identify a set of structural features for PE files which is computable in realtime, (2) use an efficient preprocessor for removing redundancy in the features' set, and (3) select an efficient data mining algorithm for final classification. Consequently, our proposed framework consists of three modules: the feature extraction module, the feature selection/preprocessing module, and the detection module.

We have evaluated our proposed detection framework on two independently collected malware datasets with different statistics. The first malware dataset is the VX Heavens Virus collection consisting of more than ten thousand malicious PE files [27]. The second malware dataset is the Malfease dataset, which contains more than five thousand malicious PE files [21]. We also collected more than one thousand benign PE files from our virology lab, which we use in conjunction with both malware datasets in our study. The results of our experiments

---

<sup>1</sup> We define a technique as *realtime deployable* if it has three properties: (1) a *tp* rate (or true positive rate) of approximately 1, (2) an *fp* rate (or false positive rate) of approximately 0, and (3) the file scanning time is comparable to existing COTS AV.

show that our PE-miner framework achieves more than 99% detection rate with less than 0.5% false alarm rate for distinguishing between benign and malicious executables. Further, our framework takes on the average only 0.244 seconds to scan a given PE file. Therefore, we can conclude that PE-Miner is *realtime deployable*, and consequently it can be easily integrated into existing COTS AV products. PE-Miner framework can also categorize the malicious executables as a function of their payload. This analysis is of great value for system administrators and malware forensic experts. An interested reader can find details in the accompanying technical report [23].

We have also compared PE-Miner with other promising malware detection schemes proposed by Perdisci et al. [18], Schultz et al. [22], and Kolter et al. [11]. These techniques use some variation of  $n$ -gram analysis for malware detection. PE-Miner provides better detection accuracy<sup>2</sup> with significantly smaller processing overheads compared with these approaches. We believe that the superior performance of PE-Miner is attributable to a rich set of novel PE format specific structural features, which provides relevant information for better detection accuracy [10]. In comparison,  $n$ -gram based techniques are more suitable for classification of loosely structured data; therefore, they fail to exploit format specific structural information of a PE file. As a result, they provide lower detection rates and have higher processing overheads as compared to PE-Miner. Our experiments also demonstrate that the detection mechanism of PE-Miner does not show any significant bias towards packed/non-packed PE files. Finally, we investigate the robustness of PE-Miner against “crafty” attacks which are specifically designed to evade detection mechanism of PE-Miner. Our results show that PE-Miner is resilient to majority of such evasion attacks.

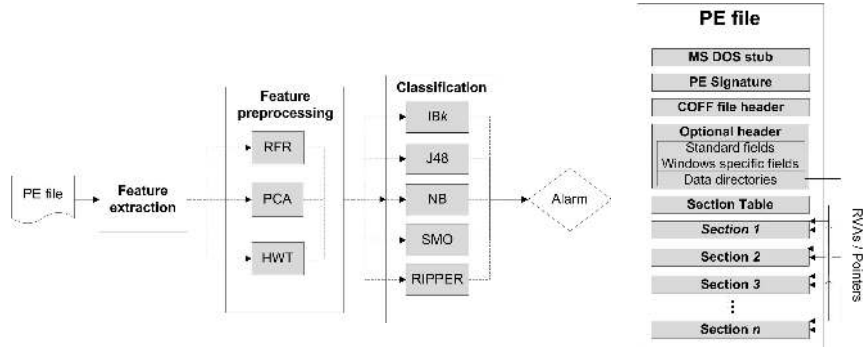
## 2 PE-Miner Framework

In this section, we discuss our proposed PE-Miner framework. We set the following strict requirements on our PE-Miner framework to ensure that our research is enacted with a product development cycle that has a short time-to-market:

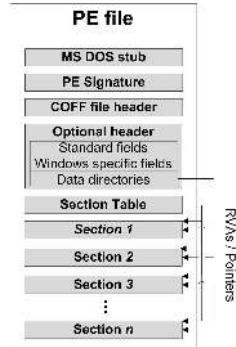
- It must be a pure non-signature based framework with an ability to detect zero-day malicious PE files.
- It must be *realtime deployable*. To this end, we say that it should have more than 99%  $tp$  rate and less than 1%  $fp$  rate. We argue that it is still a challenge for non-signature based techniques to achieve these true and false positive rates. Moreover, its time to scan a PE file must be comparable to those of existing COTS AV products.

---

<sup>2</sup> Throughout this text, the terms *detection accuracy* and *Area Under ROC Curve (AUC)* are used interchangeably. ROC curves are extensively used in machine learning and data mining to depict the tradeoff between the true positive rate and false positive rate of a classifier. The AUC ( $0 \leq \text{AUC} \leq 1$ ) is used as a yardstick to determine the detection accuracy from ROC curve. Higher values of AUC mean high  $tp$  rate and low  $fp$  rate [28]. At  $\text{AUC} = 1$ ,  $tp$  rate = 1 and  $fp$  rate = 0.



**Fig. 1.** The architecture of our PE-Miner framework



**Fig. 2.** The PE file format

- Its design must be modular that allows for the plug-n-play design philosophy. This feature will be useful in customizing the detection framework to specific requirements, such as porting it to the file formats used by other operating systems.

We have evolved the final modular architecture of our PE-Miner framework in a question oriented engineering fashion. In our research, we systematically raised following relevant questions, analyzed their potential solutions, and finally selected the best one through extensive empirical studies.

1. Which PE format specific features can be statically extracted from PE files to distinguish between benign and malicious files? Moreover, are the format specific features better than the existing  $n$ -grams or string-based features in terms of detection accuracy and efficiency?
2. Do we need to deploy preprocessors on the features' set? If yes then which preprocessors are best suited for the raw features' set?
3. Which are the best back-end classification algorithms in terms of detection accuracy and processing overheads.

Our PE-Miner framework consists of three main modules inline with the above-mentioned vision: (1) feature extraction, (2) feature preprocessing, and (3) classification (see Figure 1). We now discuss each module separately.

## 2.1 Feature Extraction

Let us revisit the PE file format [12] before we start discussing the structural features used in our features' set. A PE file consists of a PE file header, a section table (section headers) followed by the sections' data. The PE file header consists of a MS DOS stub, a PE file signature, a COFF (Common Object File Format) header, and an optional header. It contains important information about a file such as the number of sections, the size of the stack and the heap, etc. The section table contains important information about the sections that follow it,

Table 1. List of the features extracted from PE files

Feature Description	Type	Quantity
DLLs referred	binary	73
COFF file header	integer	7
Optional header – standard fields	integer	9
Optional header – Windows specific fields	integer	22
Optional header – data directories	integer	30
.text section – header fields	integer	9
.data section – header fields	integer	9
.rsrc section – header fields	integer	9
Resource directory table & resources	integer	21
<b>Total</b>		189

Table 2. Mean values of the extracted features. The bold values in every row highlight interesting outliers.

Dataset Name of Feature	VX Heavens										Malfease -
	Benign	Backdoor + Sniffer	Constructor + Virtool	DoS + Nuker	Flooder	Exploit + Hacktool	Worm	Trojan	Virus		
WSOCK32.DLL	0.037	<b>0.503</b>	0.038	<b>0.188</b>	<b>0.353</b>	<b>0.261</b>	<b>0.562</b>	<b>0.242</b>	0.053	0.065	
WININET.DLL	0.073	<b>0.132</b>	0.009	0.013	0.04	<b>0.141</b>	0.004	<b>0.103</b>	0.019	0.086	
# Symbols	<b>430.2</b>	<b>2.0E6</b>	14.7	59.4	25.8	<b>3.5E6</b>	38.8	<b>4.1E6</b>	<b>1.0E6</b>	<b>2.7E7</b>	
Maj Linker Ver	<b>4.7</b>	14.4	11.2	14.1	12.1	12.3	18.7	12.2	19.3	6.5	
Init Data Size (E5)	<b>4.4</b>	1.1	0.5	0.4	0.8	0.7	0.4	0.4	0.1	0.6	
Maj Img Ver	<b>163.1</b>	1.6	6.3	0.4	0.6	11.2	0.3	6.0	53.6	0.2	
DLL Char	<b>5.8x10<sup>3</sup></b>	0.0	0.0	0.0	0.0	24.9	0.0	3.1	230.8	18.7	
Exp Tbl Size (E2)	<b>13.7</b>	2.4	1.7	<b>14.1</b>	5.0	0.3	1.2	2.1	0.9	0.05	
Imp Tbl Size (E2)	5.8	<b>19.2</b>	6.1	7.9	<b>20.8</b>	7.1	<b>23.4</b>	<b>10.3</b>	6.2	4.7	
Rsrc Tbl Size (E4)	<b>32.6</b>	5.5	1.5	1.4	6.2	1.0	2.6	2.2	0.5	5.9	
Except Tbl Size	12.0	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	<b>3.5</b>	
.data Raw Size (E3)	<b>25.2</b>	8.4	5.6	6.3	6.0	7.9	6.1	5.5	6.7	22.1	
# Cursors	<b>14.5</b>	6.4	6.7	7.4	6.1	5.9	5.8	6.0	3.0	6.8	
# Bitmaps	<b>12.6</b>	1.2	0.0	1.0	0.6	0.7	1.2	1.4	2.4	0.5	
# Icons	<b>17.6</b>	2.5	1.9	2.7	2.0	2.1	1.8	1.9	4.5	2.2	
# Dialogs	<b>10.9</b>	3.2	1.5	3.2	1.5	2.0	1.9	1.7	2.2	2.3	
# Group Cursors	<b>11.6</b>	6.0	6.6	7.2	5.8	5.8	5.4	5.7	2.7	6.7	
# Group Icons	<b>4.1</b>	1.0	0.7	1.0	0.8	0.7	0.5	0.7	1.5	0.6	

such as their name, offset and size. These sections contain the actual data such as code, initialized data, exports, imports and resources [12], [15].

Figure 2 shows an overview of the PE file format [12], [15]. It is important to note that the section table contains Relative Virtual Addresses (RVAs) and the pointers to the start of every section. On the other hand, the data directories in an optional header contain references to various tables (such as import, export, resource, etc.) present in different sections. These references, if appropriately analyzed, can provide useful information.

We believe that this structural information about a PE file should be leveraged to extract features that have the potential to achieve high detection accuracy. Using this principle, we statically extract a set of large number of features from a given PE file<sup>3</sup>. These features are summarized in Table 1. In the discussion below, we first intuitively argue about the features that have the potential to distinguish between benign and malicious files. We then show interesting observations derived from the executable datasets used in our empirical studies.

**DLLs referred.** The list of DLLs referred in an executable effectively provides an overview of its functionality. For example, if an executable calls `WINSOCKET.DLL`

<sup>3</sup> A well-known Microsoft Visual C++ utility, called `dumpbin`, dumps the relevant information which is present inside a given PE file [4]. Another freely available utility, called `pedump`, also does the required task [20].

or `WSOCK.DLL` then it is expected to perform network related activities. However, there can be exceptions to this assumption as well. In [22], Schultz et al. have used the conjunction of DLL names, with a similar functionality, as binary features. The results of their experiments show that this feature helps to attain reasonable detection accuracy. However, our pilot experimental studies have revealed that using them as individual binary features can reveal more information, and hence can be more helpful in detecting malicious PE files. In this study, we have used 73 core functionality DLLs as features. Their list and functionality is detailed in [23]. Table 2 shows the mean feature values for the two DLLs<sup>4</sup>. Interestingly, `WSOCK32.DLL` and `WININET.DLL` are used by the majority of backdoors, nukers, flooders, hacktools, worms, and trojans to access the resources on the network and the Internet. Therefore, the applications misusing these DLLs might provide a strong indication of a possible covert network activity.

**COFF file header.** The COFF file header contains important information such as the type of the machine for which the file is intended, the nature of the file (DLL, EXE, or OBJ etc.), the number of sections, and the number of symbols. It is interesting to note in Table 2 that a reasonable number of symbols are present in benign executables. The malicious executables, however, either contain too many or too few symbols.

**Optional header: standard fields.** The interesting information in the standard fields of the optional header include the linker version used to create an executable, the size of the code, the size of the initialized data, the size of the uninitialized data, and the address of the entry point. Table 2 shows that the values of major linker version and the size of the initialized data have a significant difference in the benign and malicious executables. The size of the initialized data in benign executables is usually significantly higher compared to those of the malicious executables.

**Optional header: Windows specific fields.** The Windows specific fields of the optional header include information about the operating system version, the image version, the checksum, the size of the stack and the heap. It can be seen in Table 2 that the values of fields such as the major image version and the DLL characteristics are usually set to zero in the malicious executables. In comparison, their values are significantly higher in the benign executables.

**Optional header: data directories.** The data directories of the optional header provide pointers to the actual data present in the sections following it. It includes the information about export, import, resource, exception, debug, certificate, and base relocation tables. Therefore, it effectively provides a summary of the contents of an executable. Table 2 highlights that the size of the export table is higher for the benign executables and nukers as compared to those of other malicious executables. Another interesting observation in Table 2 is that the backdoors, flooders, worms and trojans mostly have a bigger import table size. It can be intuitively argued that they usually import network functionalities which increase the size of their import table. The size of the resource table, on the other hand, is higher for the benign executables as compared to those of

<sup>4</sup> The details of the datasets and their categorization are available in Section 3.

the malicious executables. The exception table is mostly absent in the malicious executables.

**Section headers.** The section headers provide important characteristics of a section such as its address, size, number of relocations, and line numbers. In this study, we have only considered text, data and resource sections because they are commonly present in the executables. Note that the size of the data section (if present) is relatively larger for the benign executables.

**Resource directory table & resources.** The resource directory table provides an overview of the resources that are present in the resource section of an executable file. We consider the actual count of various types of resources that are present in the resource section of an executable file. The typical examples of resources include cursors, bitmaps, icons, menus, dialogs, fonts, group cursors, and user defined resources. Intuitively and as shown in Table 2, the number of these resources is relatively higher for the benign executables.

## 2.2 Feature Selection/Preprocessing

We have now identified our features' set that consists of a number of statically computable features – 189 to be precise – based on the structural information of the PE files. It is possible that some of the features might not convey useful information in a particular scenario. Therefore, it makes sense to remove or combine them with other similar features to reduce the dimensionality of our input feature space. Moreover, this preprocessing on the raw extracted features' set also reduces the processing overheads in training and testing of classifiers, and can possibly also improve the detection accuracy of classifiers. In this study, we have used three well-known features' selection/preprocessing filters. We provide their short descriptions in the following text. More details can be found in [29].

**Redundant Feature Removal (RFR).** We apply this filter to remove those features that do not vary at all or show significantly large variation i.e. they have approximately uniform-random behavior. Consequently, this filter removes all features that have either constant values or show a variance above a threshold or both.

**Principal Component Analysis (PCA).** The Principal Component Analysis (PCA) is a well-known filter for dimensionality reduction. It is especially useful when the input data has high dimensionality – sometimes referred to as *curse of dimensionality*. This dimensionality reduction can possibly improve the quality of an analysis on a given data if the dataset consists of highly correlated or redundant features. However, this dimensionality reduction may result in information loss (i.e. reduction in data variance) as well. One has to carefully choose the appropriate balance for this tradeoff. We apply PCA filter to remove/combine correlated features for dimensionality reduction.

**Haar Wavelet Transform (HWT).** The principle of this technique is that the most relevant information is stored with the highest coefficients at each order of a transform. The lower order coefficients can be ignored to get only the most relevant information. The wavelet transform has also been used for dimensionality reduction. The wavelet transform technique has been extensively



Table 3. Statistics of the data used in this study.

Dataset	VX Heavens									Malfease
	Benign	Backdoor + Sniffer	Constructor + Virtool	DoS + Nuker	Flooder	Exploit + Hacktool	Worm	Trojan	Virus	
Quantity	1,447	3,455	367	267	358	243	1,483	3,114	1,052	5,586
Avg. Size (KB)	1,263	270	234	176	298	156	72	136	50	285
Min. Size (KB)	4	1	4	3	6	4	2	1	2	1
Max. Size (KB)	104,588	9,277	5,832	1,301	14,692	1,924	2,733	4,014	1,332	5,746
UPX	17	786	79	15	32	43	353	622	48	470
ASPack	2	432	21	16	25	15	66	371	10	187
Misc. Packed	372	325	47	31	58	38	471	170	71	1,909
Borland C/C++	15	56	8	15	10	6	13	63	18	11
Borland Delphi	13	589	13	65	64	8	76	379	71	342
Visual Basic	4	719	106	39	126	38	210	674	119	809
Visual C++	526	333	19	51	29	59	89	619	96	351
Visual C#	56	0	0	0	1	0	5	1	6	1
Misc. Other	9	49	9	2	3	2	4	15	7	5
Non-packed (%)	43.1	50.5	42.2	64.4	65.1	46.5	26.8	56.2	30.1	27.2
Packed (%)	27.0	44.7	40.1	23.2	32.1	39.5	60.0	37.4	12.3	46.6
Not Found (%)	29.9	4.8	17.7	12.4	2.8	14.0	13.2	6.4	57.6	26.2

used in the image compression but is never evaluated in the malware detection domain. The Haar wavelet is one of the simplest wavelets and is known to provide reasonable accuracy. The application of Haar wavelet transform requires input data to be normalized. Therefore, we have passed the data through a *normalize* filter before applying HWT.

### 2.3 Classification

Once the dimensionality of the input features' set is reduced by applying one of the above-mentioned preprocessing filters, it is given as an input to the well-known data mining algorithms for classification. In this study we have used five classifiers: (1) instance based learner (IBk), (2) decision tree (J48), (3) Naïve Bayes (NB), (4) inductive rule learner (RIPPER), and (5) support vector machines using sequential minimal optimization (SMO). An interested reader can find their details in the accompanying technical report [23].

## 3 Datasets

In this section, we present an overview of the datasets used in our study. We have collected 1,447 benign PE files from the local network of our virology lab. The collection contains executables such as Packet CAPture (PCAP) file parsers compiled by MS Visual Studio 6.0, compressed installation executables, and MS Windows XP/Vista applications' executables. The diversity of the benign files is also evident from their sizes, which range from a minimum of 4 KB to a maximum of 104,588 KB (see Table 3).

Moreover, we have used two malware collections in our study. First is the VX Heavens Virus Collection, which is *labeled* and is publicly available for free download [27]. We only consider PE files to maintain focus. Our filtered dataset

contains 10,339 malicious PE files. The second dataset is the Malfease malware dataset [21], which consists of 5,586 *unlabeled* malicious PE files.

In order to conduct a comprehensive study, we further categorize the malicious PE files as a function of their payload<sup>5</sup>. The malicious executables are subdivided into eight major categories such as *virus*, *trojan*, *worm*, etc [7]. Moreover, we have combined some categories that have similar functionality. For example, we have combined *constructor* and *virtool* to create a single *constructor + virtool* category. This unification increases the number of malware samples per category. Brief introductions of every malware category are provided in the accompanying technical report [23].

Table 3 provides the detailed statistics of the malware used in our study. It can be noted that the average size of the malicious executables is smaller than that of the benign executables. Further, some executables used in our study are encrypted and/or compressed (packed). The detailed statistics about packing are also tabulated in Table 3. We use PEiD [16] and Protection ID for detecting packed executables [19]<sup>6</sup>.

Our analysis shows that VX Heavens Virus collection contains 40.1% packed and 47.2% non-packed PE files. However, approximately 12.7% malicious PE files cannot be classified as either packed or non-packed by PEiD and Protection ID. The Malfease collection contains 46.6% packed and 27.2% non-packed malicious PE files. Similarly, 26.2% malicious PE files cannot be classified as packed or non-packed. Therefore, we can say that packed/non-packed malware distribution in the VX Heavens virus collection is relatively more balanced than the Malfease dataset. In our collection of benign files, 43.1% are packed and 27.0% are non-packed PE files respectively. Similarly, 29.9% benign files are not detected by PEiD and Protection ID. An interesting observation is that the benign PE files are mostly packed using nonstandard and custom developed packers. We speculate that a significant portion of the packed executables are not classified as packed because the signatures of their respective packers are not present in the database of PEiD or Protection ID. *Note that we do not manually unpack any PE file prior to the processing of our PE-Miner.*

## 4 Related Work

We now briefly describe the most relevant non-signature based malware detection techniques. These techniques are proposed by Perdisci et al. [18], Schultz et al. [22] and Kolter et al. [11]. We briefly summarize their working principles in the following paragraphs but an interested reader can find their detailed description in [23].

In [18], the authors proposed McBoost that uses two classifiers – C1 and C2 – for classification of non-packed and packed PE files respectively. A custom

<sup>5</sup> Since the Malfease malware collection is unlabeled; therefore, it is not possible to divide it into different malware categories.

<sup>6</sup> We acknowledge the fact that PEiD and Protection ID are signature based packer detectors and can have significant false negatives.

developed unpacker is used to extract the hidden code from the packed PE files and the output of the unpacker is given as an input to the C2 classifier. Unfortunately, we could not obtain its source code or binary due to licensing related problems. Furthermore, its implementation is not within the scope of our current work. Consequently, we only evaluate the C1 module of McBoost which works only for non-packed PE files. Therefore, we acknowledge that our McBoost results should be considered only preliminary.

In [22], Schultz et al. have proposed three independent techniques for detecting malicious PE files. The first technique, uses the information about DLLs, function calls and their invocation counts. However, the authors did not provide enough information about the used DLLs and function names; therefore, it is not possible for us to implement it. But we have implemented the second approach (titled *strings*) which uses strings as binary features i.e. present or absent. The third technique uses two byte words as binary features. This technique is later improved in a seminal work by Kolter et al. [11] which uses 4-grams as binary features. Therefore, we include the technique of Kolter et al. (titled *KM*) in our comparative evaluation.

## 5 Experimental Results

We have compared our PE-Miner framework with recently proposed promising techniques by Perdisci et al. [18], Schultz et al. [22], and Kolter et al. [11]. We have used the standard 10 fold cross-validation process in our experiments, i.e., the dataset is randomly divided into 10 smaller subsets, where 9 subsets are used for training and 1 subset is used for testing. The process is repeated 10 times for every combination. This methodology helps in systematically evaluating the effectiveness of our approach to detect previously unknown (i.e. zero-day) malicious PE files. The ROC curves are generated by varying the threshold on output class probability [5], [28]. The AUC is used as a yardstick to determine the detection accuracy of each approach. We have done the experiments on an Intel Pentium Core 2 Duo 2.19 GHz processor with 2 GB RAM. The Microsoft Windows XP SP2 is installed on this machine.

### 5.1 Malicious PE File Detection

In our first experimental study, we attempt to distinguish between benign and malicious PE files. To get better insights, we have done independent experiments with benign and each of the eight types of the malicious executables. The five data mining algorithms, namely IBk, J48, NB, RIPPER, and SMO, are deployed on top of each approach (namely PE-Miner with RFR, PE-Miner with PCA, PE-Miner with HWT, McBoost (C1 only) by Perdisci et al. [18], strings approach by Schultz et al. [22], and KM by Kolter et al. [11]). This results in a total of 270 experimental runs each with 10-fold cross validation. We tabulate our results for this study in Table 4 and now answer different questions that we raised in Section 2 in a chronological fashion.

Table 4. AUCs for detecting the malicious executables. The bold entries in each column represent the best results.

Dataset	VX Heavens									Malfease
Malware	Backdoor + Sniffer	Constructor + Virtool	DoS + Nuker	Flooder	Exploit + Hacktool	Worm	Trojan	Virus	Average	-
PE-Miner — RFR										
IBK	0.992	0.996	0.995	0.994	0.998	0.979	0.984	0.994	<b>0.992</b>	0.986
J48	0.993	<b>0.998</b>	0.987	0.993	0.999	0.979	<b>0.992</b>	0.993	<b>0.992</b>	0.979
NB	0.971	0.978	0.966	0.973	0.987	0.972	0.974	0.986	0.976	0.976
RIPPER	<b>0.996</b>	0.996	0.977	0.981	0.999	<b>0.988</b>	0.988	0.996	0.990	0.985
SMO	0.991	0.990	0.991	0.993	0.997	0.975	0.978	0.992	0.988	0.963
PE-Miner — PCA										
IBK	0.989	0.996	0.994	0.995	0.998	0.976	0.984	0.993	0.991	0.984
J48	0.980	0.966	0.929	0.960	0.987	0.936	0.951	0.985	0.962	0.945
NB	0.961	0.990	0.993	0.996	0.996	0.964	0.956	0.990	0.981	0.898
RIPPER	0.982	0.978	<b>0.996</b>	0.974	0.977	0.949	0.968	0.987	0.976	0.952
SMO	0.990	0.992	0.989	0.995	0.995	0.958	0.965	0.992	0.985	0.954
PE-Miner — HWT										
IBK	0.991	0.996	<b>0.996</b>	<b>0.998</b>	<b>1.000</b>	0.978	0.985	0.995	<b>0.992</b>	0.986
J48	0.995	0.997	0.993	0.988	0.997	0.978	0.991	0.999	<b>0.992</b>	0.977
NB	0.989	0.982	0.983	0.987	0.990	0.978	0.972	0.990	0.984	0.960
RIPPER	0.994	0.997	0.982	0.990	0.997	0.983	0.990	<b>1.000</b>	<b>0.992</b>	<b>0.987</b>
SMO	0.990	0.995	0.991	0.996	<b>1.000</b>	0.972	0.973	0.994	0.989	0.964
McBoost — C1 only										
IBK	0.941	0.935	0.875	0.960	0.832	0.938	0.930	0.914	0.916	0.949
J48	0.866	0.895	0.809	0.893	0.731	0.906	0.902	0.882	0.860	0.860
NB	0.831	0.924	0.723	0.889	0.795	0.873	0.886	0.844	0.846	0.817
RIPPER	0.833	0.888	0.744	0.918	0.660	0.866	0.838	0.844	0.824	0.860
SMO	0.802	0.887	0.759	0.910	0.678	0.854	0.805	0.827	0.815	0.835
Strings										
IBK	0.949	0.860	0.902	0.980	0.925	0.928	0.863	0.952	0.920	0.944
J48	0.913	0.834	0.862	0.695	0.871	0.908	0.836	0.938	0.857	0.929
NB	0.920	0.830	0.882	0.726	0.886	0.901	0.828	0.905	0.860	0.930
RIPPER	0.843	0.797	0.714	0.578	0.712	0.892	0.743	0.929	0.776	0.927
SMO	0.855	0.817	0.705	0.775	0.583	0.871	0.756	0.883	0.781	0.933
KM										
IBK	0.984	0.934	0.983	0.971	0.983	0.987	0.979	0.986	0.976	0.980
J48	0.953	0.940	0.916	0.907	0.916	0.957	0.951	0.953	0.937	0.952
NB	0.943	0.959	0.961	0.952	0.961	0.968	0.954	0.954	0.957	0.961
RIPPER	0.951	0.944	0.924	0.921	0.924	0.964	0.948	0.948	0.941	0.971
SMO	0.949	0.946	0.952	0.927	0.952	0.961	0.940	0.938	0.946	0.960

**Which features' set is the best?** Table 4 tabulates the AUCs for PE-Miner using three different preprocessing filters (RFR, PCA and HWT), McBoost, strings and KM [11]. A macro level scan through the table clearly shows the supremacy of PE-Miner based approaches with AUCs more than 0.99 for most of the malware types and even approaching 1.00 for some malware types. For PE-Miner, RFR and HWT preprocessing lead to the best average results with more than 0.99 AUC.

The strings approach gives the worst detection accuracy. The KM approach is better than the strings approach but inferior to our PE-Miner. This is expected because the string features are not stable as compiling a given piece of code by using different compilers leads to different sets of strings. Our analysis shows that KM approach is more resilient to variation in the string sets because it uses a combination of string and non-string features. The results obtained for KM approach (AUC= 0.95) are also consistent with the results reported in [11]. The C1 module of McBoost also provides relatively inferior detection accuracies which are as low as 0.66 for exploit+hacktool category. It is important to note that the C1 module of McBoost is functionally similar to the techniques proposed by Schultz et al. and Kolter et al. The only significant difference is that C1 operates only on the code sections of the non-packed PE files whereas the other techniques operate on complete files.

**Table 5.** The processing overheads (in seconds/file) of different feature selection, extraction and preprocessing schemes.

	PE-Miner			McBoost	Strings	KM
	(RFR)	(PCA)	(HWT)			
Selection	-	-	-	2.839	5.289	31.499
Extraction	0.228	0.228	0.228	0.198	0.130	0.220
Preprocessing	0.007	0.009	0.012	-	-	-
Total	<b>0.235</b>	0.237	0.240	3.037	5.419	31.719

It is important to emphasize that both strings and KM approaches incur large overheads in the feature selection process (see Table 5<sup>7</sup>). Kolter et al. have confirmed that their implementation of information gain calculation for feature selection took almost a day for every run. To make our implementation of  $n$ -grams more efficient, we use `hash_map` STL containers in the Visual C++ [8]. Our experiments show that the feature selection process in KM still takes more than 31 seconds per file even with our optimized implementation. The optimized strings approach takes, on the average, more than 5 seconds per file for feature selection. The optimized McBoost (C1 only) approach takes an average of more than 2 seconds per file for feature selection<sup>8</sup>. These approaches have processing overheads because the time to calculate information gain increases exponentially with the number of unique  $n$ -grams (or strings). On the other hand, PE-Miner does not suffer from such serious bottlenecks. The application of RFR, PCA, or HWT filters takes only about a hundredth of a second.

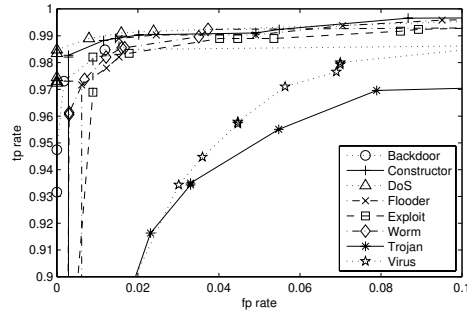
**Which classification algorithm is the best?** We can conclude from Table 4 that J48 outperforms the rest of the data mining classifiers in terms of the detection accuracy in most of the cases. Moreover, Table 6 shows that J48 has one of the smallest processing overheads both in training and testing. RIPPER and IBk closely follow the detection accuracy of J48. However, they are infeasible for realtime deployment because of the high processing overheads in the training and the testing phases respectively. The processing overheads of training RIPPER are the highest among all classifiers. In comparison, IBk does not require a training phase but its processing overheads in the testing phase are the highest. Further, Naïve Bayes gives the worst detection accuracy because it assumes independence among input features. Intuitively speaking, this assumption does not hold for the features' sets used in our study. Note that Naïve Bayes has very small learning and testing overheads (see Table 6<sup>9</sup>).

**Which malware category is the most challenging to detect?** An overview of Table 4 suggests that the most challenging malware categories are worms and trojans. The average AUC values of the compared techniques for worms and trojans are approximately 0.95. The poor detection accuracy is attributed to the fact that the trojans are inherently designed to appear similar to the benign

<sup>7</sup> The results in Table 5 are averaged over 100 runs.

<sup>8</sup> Note that the complete McBoost system also uses unpacker for extraction of hidden code. This process is time consuming as reported by the authors in [18].

<sup>9</sup> The results in Table 6 are averaged over 100 runs.



**Fig. 3.** The magnified ROC plots for detecting the malicious executables using PE-Miner utilizing J48 preprocessed with RFR filter.

**Table 6.** The processing overheads (in seconds/file) of different features and classification algorithms.

	IBK	J48	NB	RIPPER	SMO	IBK	J48	NB	RIPPER	SMO
<b>Training</b>					<b>Testing</b>					
PE-Miner (RFR)	-	0.008	0.001	0.269	0.199	0.032	0.001	0.002	0.002	0.002
PE-Miner (PCA)	-	0.007	0.001	0.264	0.179	0.035	0.001	0.001	0.001	0.002
PE-Miner (HWT)	-	0.007	0.001	0.252	0.147	0.032	0.001	0.002	0.001	0.002
McBoost	-	0.021	0.004	1.305	1.122	0.218	0.010	0.007	0.005	0.022
Strings	-	0.009	0.002	0.799	0.838	0.163	0.003	0.003	0.002	0.003
KM	-	0.024	0.004	1.510	1.018	0.254	0.018	0.007	0.005	0.020

**Table 7.** Realtime deployable analysis of the best techniques

Technique	Classifier	AUC	Scan Time (sec/file)	Is Realtime Deployable?
PE-Miner (RFR)	J48	0.991	0.244	“Yes”
McBoost	IBk	0.926	3.255	No
Strings	IBk	0.927	5.582	No
KM	IBk	0.977	31.973	No
AVG Free 8.0 [1]	-	-	0.159	-
Panda 7.01 [14]	-	-	0.131	-

executables. Therefore, it is a difficult challenge to distinguish between trojans and benign PE files. Our PE-Miner still achieves on the average 0.98 AUC for worms and trojans which is quite reasonable. Figure 3 shows that for other malware categories, PE-Miner (with RFR preprocessor) has AUCs more than 0.99.

## 5.2 Miscellaneous Discussions

We conclude our comparative study with an answer to an important issue: *which of the compared techniques meet the criterion of being realtime deployable?* (see Section 2). We tabulate the AUC and the scan time of the best techniques in Table 7. Moreover, we also show the scan time of two well-known COTS AV products for doing the *realtime deployable* analysis of different non-signature based techniques. It is clear that PE-Miner (RFR) with J48 classifier is the only non-signature based technique that satisfies the criterion of being *realtime deployable*. One might argue that PE-Miner framework provides only a small

**Table 8.** Portion of the developed decision trees for distinguishing between benign and backdoor+sniffer

```

NumMessageTable <= 0
|   SizeLoadConfigTable <= 0
|   |   TimeDateStamp <= 1000000000
|   |   |   NumCursor <= 1
|   |   |   |   NumAccelerators <= 0
|   |   |   |   |   NumBitmap <= 0: malicious
|   |   |   |   |   NumBitmap > 0: benign
|   |   |   |   |   NumAccelerators > 0:malicious
|   |   |   |   NumCursor > 1:malicious

```

improvement in detection accuracy over the KM approach. *But then KM has the worst scan time of 31.97 seconds per file (see Table 7).* It is very important to interpret the results in Table 7 from a security expert’s perspective. For example, if a malware detector scans ten thousand files with an AUC of 0.97, it will not detect approximately 300 malicious files. In comparison, a detector with an AUC of 0.99 will miss only 100 files, which is a 66.6% improvement in the number of missed files [2]. Therefore, we argue that from a security expert’s perspective, even a small improvement in the detection accuracy is significant in the limiting case when the detection accuracy approaches to 1.00.

An additional benefit of PE-Miner is that it provides insights about the learning models of different classifiers that can be of great value to malware forensic experts. We show a partial subtree of J48 for categorizing benign and malicious PE files in Table 8. The message tables mostly do not exist in the backdoor+sniffer categories. The TimeDateStamp is usually obfuscated in the malicious executables. The number of resources are generally smaller in malicious PE files, whereas the benign files tend to have larger number of resources such as menus, icons, and user defined resources. Similar insights are also provided by the rules developed in the training phase of RIPPER.

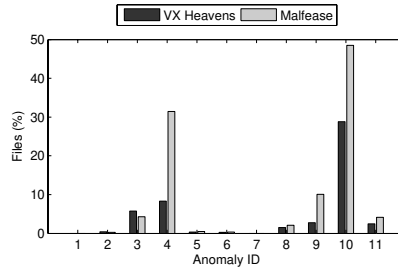
In [9], the authors have pointed out several difficulties in parsing PE files. In our experiments, we have also observed various anomalies in parsing the structure of malicious PE files. Table 9 contains the list of anomalies which we have observed in parsing malicious PE files. A significant proportion of malicious PE files have anomalous structure which can crash a naïve PE file parser. Figure 4 provides the statistics of anomalies which we have observed in parsing malicious PE files of VX Heavens and Malfease collections. To this end, we have developed a set of heuristics which successfully handle the above-mentioned anomalies.

## 6 Robustness and Reliability Analysis of PE-Miner

We have now established the fact that *PE-Miner is a realtime deployable scheme for zero-day malware detection.* A careful reader might ask whether the statement still holds if the “ground truth” is now changed as: (1) we cannot trust the classification of signature based packer detectors PEiD and Protection ID, and (2) a “crafty” attacker can forge the features of malicious files with those of

**Table 9.** List of the anomalies observed in parsing malicious PE files

ID	Description
1	Large number of sections
2	SizeOfHeader field is unaligned
3	Overlapping DoS and PE headers
4	Large virtual size in a section
5	Large raw data size in a section
6	Zero/Non-zero pair in data directory table
7	Large pointer in data directory entry
8	Size of section is too large
9	Section name garbled (non printable characters)
10	There is an unknown overlay region
11	Out of file pointer

**Fig. 4.** Statistics of anomalies observed in parsing malicious PE files

benign files to evade detection. In this section, we do a stress and regression testing of PE-Miner to analyze robustness of its features and its resilience to potential evasive techniques.

### 6.1 Robustness Analysis of Extracted Features

It is a well-known fact that signature based packer detector PEiD, which we are using to distinguish between packed and non-packed executables, has approximately 30% false negative rate [17]. In order to convince ourselves that our extracted features are actually “robust”, we evaluate PE-Miner in four scenarios: (1) training PE-Miner on 70% non-packed PE files and 30% packed PE files and testing on the remaining 70% packed PE files, (2) training PE-Miner on non-packed PE files only and testing on packed PE files, (3) training PE-Miner on packed PE files only and then testing on non-packed PE files, and (4) testing PE-Miner on a “difficult” dataset that consists of packed benign and non-packed malicious PE files. We assert that the scenarios (2) and (3) – although unrealistic – still provide valuable insight into the extent of bias, that PE-Miner might have, towards detection of packed/non-packed executables.

We want to emphasize an important point that there is no confusion about “ground truth” for packed executables in above-mentioned four scenarios because a packer only detects a file as “packed” if it has its signature in its database. The confusion about “ground truth”, however, stems in the fact that a reasonable proportion of packed PE files could be misclassified as non-packed because of



Table 10. An analysis of robustness of extracted features of PE-Miner (RFR) in different scenarios

Dataset	VX Heavens									Malfease
Malware	Backdoor + Sniffer	Constructor + Virtool	DoS + Nuker	Flooder	Exploit + Hacktool	Worm	Trojan	Virus	Average	-
<b>Scenario 1: Detection of packed benign and malicious PE files</b>										
IBK	0.999	1.000	1.000	0.999	0.999	0.998	0.999	0.999	0.999	0.812
J48	0.996	1.000	1.000	0.999	0.999	0.998	0.993	0.999	<b>0.998</b>	<b>0.991</b>
NB	0.971	0.988	0.963	0.955	0.996	0.980	0.978	0.987	0.977	0.934
RIPPER	0.997	0.996	0.999	0.990	0.993	0.985	0.858	0.998	0.977	0.988
SMO	0.985	0.998	1.000	0.996	0.994	0.994	0.985	0.998	0.994	0.706
<b>Scenario 2: Training using non-packed executables only and testing using packed executables</b>										
IBK	0.986	0.965	0.912	0.963	0.998	0.993	0.850	0.989	0.957	0.917
J48	0.982	0.999	0.998	0.937	0.999	0.963	0.857	0.954	<b>0.961</b>	<b>0.968</b>
NB	0.927	0.899	0.842	0.809	0.966	0.911	0.857	0.965	0.897	0.780
RIPPER	0.989	0.995	0.998	0.995	0.986	0.962	0.858	0.853	0.954	0.937
SMO	0.983	0.772	0.905	0.691	0.996	0.737	0.651	0.852	0.823	0.859
<b>Scenario 3: Training using packed executables only and testing using non-packed executables</b>										
IBK	0.975	0.965	0.964	0.878	0.793	0.982	0.911	0.904	0.921	0.855
J48	0.951	0.908	0.919	0.940	0.726	0.958	0.903	0.881	<b>0.898</b>	<b>0.903</b>
NB	0.685	0.965	0.668	0.633	0.689	0.979	0.688	0.688	0.749	0.789
RIPPER	0.979	0.938	0.967	0.972	0.747	0.768	0.840	0.867	0.885	0.904
SMO	0.977	0.941	0.877	0.882	0.536	0.983	0.835	0.904	0.867	0.849
<b>Scenario 4: Detection of packed benign and non-packed malicious PE files ("difficult" dataset)</b>										
IBK	0.999	1.000	1.000	0.999	0.998	0.998	0.994	0.998	0.998	0.992
J48	0.997	0.986	0.999	0.999	0.999	0.999	0.989	0.993	<b>0.995</b>	<b>0.996</b>
NB	0.954	0.963	0.995	0.988	0.966	0.990	0.975	0.986	0.977	0.948
RIPPER	0.998	0.984	0.998	0.993	0.986	0.999	0.992	0.996	0.993	0.948
SMO	0.989	0.996	1.000	0.997	0.996	0.997	0.984	0.992	0.994	0.945

false negative rate of PEiD. Note that the false negatives of PEiD, reported in [17], consist of two types: (1) packed PE files that are misclassified as non-packed, and (2) PE files that are unclassified. We have not included unclassified files in our dataset to remove the false negatives of the second type.

**Scenario 1: Detection of packed benign and malicious PE files.** The motivation behind the first scenario is to test if PE-Miner can distinguish between packed benign and packed malware, regardless of the type of packer. In order to ensure that our features are not influenced by the type of packing tool used to encrypt PE files, our "packed-only" dataset contains PE files (both benign and malware) packed using a variety of packers like UPX, ASPack, Armadillo, PECompact, WWPack32, Virogen Crypt 0.75, UPS-Scrambler, PEBundle and PEPack etc. Moreover, the "packed-only" dataset contains on the average 44% and 56% packed malicious and benign PE files respectively. We train PE-Miner on 70% non-packed executables and 30% packed executables and then test it on the remaining 70% packed executables. The results of our experiments for this scenario are tabulated in Table 10. We can easily conclude that PE-Miner has shown good resilience in terms of detecting accuracy once it is tested on packed benign and malicious PE files from both datasets.

**Scenarios 2 and 3: Detection of packed/non-packed malicious PE files.** In the second experiment, we train PE-Miner on non-packed benign and malicious PE files and test it on packed benign and malicious PE files. Note that this scenario is more challenging because the training dataset contains significantly less number of packed files compared with the first scenario. In the third experiment, we train PE-Miner on packed benign and malicious PE files and test on non-packed benign and malicious PE files. The results of these experiments are tabulated in Table 10. It is clear from Table 10 that the detection accuracy

of PE-Miner (RFR-J48) drops to 0.96, when it is trained on non-packed executables and tested on the packed executables. Likewise, the average detection accuracy of PE-Miner (RFR-J48) drops to 0.90 for the third scenario. Remember once we train PE-Miner on “packed only” dataset, then it gets 0% exposure to non-packed files and this explains deterioration in the detection accuracy of PE-Miner. We conclude that the detection accuracy of PE-Miner, even in these unrealistic stress testing scenarios, gracefully degrades.

**Scenario 4: Detection of packed benign and non-packed malicious PE files.** In [18], the authors report an interesting study about the ability of different schemes to detect packed/non-packed executables. They show that the detection accuracy of KM approach degrades on a “difficult” dataset consisting of packed benign and non-packed malicious PE files. According to the authors in [18], KM shows a bias towards detecting packed PE files as malware and non-packed PE files as benign. We also – in line with this strategy – tested PE-Miner on a “difficult” dataset created from both malware collections used in our study. The results are tabulated in Table 10. It is important to highlight that for these experiments PE-Miner is trained on the original datasets but is tested on the “difficult” versions of both datasets. One can conclude from the results in Table 10 that PE-Miner does not show any bias towards detecting packed executables as malicious and non-packed executables as benign.

Our experiments conclude that *the extracted features are actually “robust”, and as a result, PE-Miner does not show any significant bias towards detection of packed/non-packed executables.*

## 6.2 Reliability of PE-Miner

Now we test PE-Miner on a “crafty” malware dataset, especially designed to circumvent detection by PE-Miner. We particularly focus our attention on the *false negative rate* (or miss detection rate)<sup>10</sup> of PE-Miner when we replace features in malicious files with those of benign files. It can be argued that if adversaries exactly know our detection methodology, they might be able to design strategies that evade detection by PE-Miner. The examples of such strategies could be especially crafted packing techniques, insertion of dummy resources, obfuscation of address pointers, and other information present in headers etc.

We have conducted an empirical study to analyze the robustness of PE-Miner to such evasive techniques. To this end, we have “crafted” malware files in the datasets to contain benign-like features. Specifically, we have created seven “crafty” datasets in which for every malware file 5, 10, 30, 50, 100, 150 and 189 random features – out of 189 features – are *forged* with the respective features from a randomly chosen benign file. We now analyze the false negative rate of PE-Miner (RFR-J48) across these “crafty” datasets. The results tabulated in Table 11 highlight the robustness of PE-Miner to such crafty attacks. The false negative rate of PE-Miner stays below 1% when fifty features are simultaneously

<sup>10</sup> The false negative rate is defined by the fraction of malicious files wrongly classified as benign.

Table 11. False negative rate for detecting malicious executables with PE-Miner on the “crafty” datasets

Dataset Malware	VX Heavens									Malfease -
	Backdoor + Sniffer	Constructor + Virtool	DoS + Nuker	Flooder	Exploit + Hacktool	Worm	Trojan	Virus	Average	
# Forged Features	False negative rate									
0/189	0.001	0.000	0.000	0.000	0.004	0.000	0.004	0.007	0.002	0.001
5/189	0.002	0.000	0.000	0.000	0.004	0.000	0.004	0.007	0.002	0.001
10/189	0.002	0.000	0.000	0.000	0.004	0.011	0.004	0.014	0.004	0.004
30/189	0.002	0.003	0.000	0.012	0.023	0.011	0.011	0.014	0.009	0.004
50/189	0.002	0.003	0.000	0.012	0.023	0.016	0.011	0.014	0.010	0.004
100/189	0.096	0.003	0.000	0.012	0.023	0.050	0.445	0.176	0.101	0.004
150/189	0.658	0.003	0.000	0.583	0.795	0.611	0.558	0.221	0.429	0.426
189/189	0.996	1.000	1.000	1.000	1.000	1.000	1.000	1.000	0.999	0.998

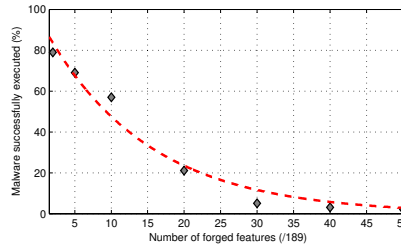


Fig. 5. Execution analysis of crafted malware files

forged. For both datasets, the average false negative rate is approximately 5% even when 100 out of 189 features are forged. This shows that a large set of features, which cover structural information of almost all portions of a PE file, used by PE-Miner make it very difficult for an attacker to evade detection – even when it manipulates majority of them at the same time.

It should be emphasized that simultaneous manipulation of all features of a PE malware file requires significant level of skill, in-depth knowledge about the structure of a PE file, and detailed understanding of our detection framework. If an attacker tries to randomly forge, using brute-force, the structural features of a PE malware file with those of a benign PE file then he/she will inevitably end up corrupting the executable image. Consequently, the file will not load successfully into memory. We have manually executed the “crafted” malicious executables. The objective is to understand that how many features a “crafty” attacker can successfully forge without ending up corrupting the executable image. The results of our experiments are shown in Figure 5. This figure proves our hypothesis that the probability of having valid PE files decreases exponentially with an increase in the number of simultaneously forged features. In fact, the successful execution probability approaches to zero as the number of simultaneously forged features approaches to 50. Referring back to Table 11, the average false negative rate of PE-Miner is less than 1% when 50 features are simultaneously forged. Therefore, we argue that it is not a cinch for an attacker to alter malicious PE files to circumvent detection by PE-Miner. However, we accept that an attacker can evade the detection capability of PE-Miner if: (1) he/she knows the exact

details of our detection framework – including the detection rules, and (2) also has the “craft” to simultaneously manipulate more than 100 structural features without corrupting the executable image.

## 7 Conclusion

In this paper we present, PE-Miner, a framework for detection of malicious PE files. PE-Miner leverages the structural information of PE files and the data mining algorithms to provide high detection accuracy with low processing overheads. Our implementation of PE-Miner completes a single-pass scan of all executables in the dataset (more than 17 thousand) in less than one hour. Therefore it meets all of our requirements mentioned in Section 2.

We believe that our PE-Miner framework can be ported to Unix and other non-Windows operating systems. To this end, we have identified similar structural features for the ELF file format in Unix and Unix-like operating systems. Our initial results are promising and show that PE-Miner framework is scalable across different operating systems. This dimension of our work will be the subject of forthcoming publications. Moreover, PE-Miner framework is also ideally suited for detecting malicious PE files on resource constrained mobile phones (running mobile variants of Windows) because of its small processing overheads. Finally, we are also doing research to develop techniques to fully remove the bias of PE-Miner in detecting packed/non-packed executables [24].

## Acknowledgments

This work is supported in part by the National ICT R&D Fund, Ministry of Information Technology, Government of Pakistan. The information, data, comments, and views detailed herein may not necessarily reflect the endorsements of views of the National ICT R&D Fund.

We are thankful to Muhammad Umer for designing experiments to collect statistics of anomalies observed in parsing malicious PE files. We also acknowledge Marcus A. Maloof and Jeremy Z. Kolter for continuous feedbacks regarding the implementation of byte sequence approach and the experimental setup. We thank Roberto Perdisci for providing implementation details of McBoost, sharing Malfease malware dataset, and the results of their custom developed unpacker. We also thank VX Heavens moderators for making a huge malware collection publicly available and sharing packing statistics of malware. We also thank Guofei Gu and Syed Ali Khayam for providing useful feedback on an initial draft of this paper.

## References

1. AVG Free Antivirus, available at <http://free.avg.com/>.
2. S. Axelsson, "The base-rate fallacy and its implications for the difficulty of intrusion detection", ACM Conference on Computer and Communications Security (CCS), pp. 1-7, Singapore, 1999.
3. J. Cheng, S.H.Y. Wong, H. Yang, S. Lu, "SmartSiren: virus detection and alert for smart-phones", International Conference on Mobile Systems, Applications and Services (MobiSys), pp. 258-271, USA, 2007.
4. DUMPBIN utility, Article ID 177429, Revision 4.0, Microsoft Help and Support, 2005.
5. T. Fawcett, "ROC Graphs: Notes and Practical Considerations for Researchers", TR HPL-2003-4, HP Labs, USA, 2004.
6. F-Secure Corporation, "F-Secure Reports Amount of Malware Grew by 100% during 2007", Press release, 2007.
7. F-Secure Virus Description Database, available at <http://www.f-secure.com/v-descs/>.
8. `hash_map`, Visual C++ Standard Library, available at <http://msdn.microsoft.com/en-us/library/6x7w9f6z.aspx>.
9. N. Hnatiw, T. Robinson, C. Sheehan, N. Suan, "PIMP MY PE: Parsing Malicious and Malformed Executables", Virus Bulletin Conference (VB), Austria, 2007.
10. K. Kendall, C. McMillan, "Practical Malware Analysis", Black Hat Conference, USA, 2007.
11. J.Z. Kolter, M.A. Maloof, "Learning to detect malicious executables in the wild", ACM International Conference on Knowledge Discovery and Data Mining (KDD), pp. 470-478, USA, 2004.
12. Microsoft Portable Executable and Common Object File Format Specification, Windows Hardware Developer Central, Updated March 2008, available at <http://www.microsoft.com/whdc/system/platform/firmware/PECOFF.mspx>.
13. J. Munro, "Antivirus Research and Detection Techniques", Antivirus Research and Detection Techniques, ExtremeTech, 2002, available at <http://www.extremetech.com/article2/0,2845,367051,00.asp>.
14. Panda Antivirus, available at <http://www.pandasecurity.com/>.
15. PE file format, Webster Technical Documentation, available at [http://webster.cs.ucr.edu/Page\\_TechDocs/pe.txt](http://webster.cs.ucr.edu/Page_TechDocs/pe.txt).
16. PEiD, available at <http://www.peid.info/>.
17. R. Perdisci, A. Lanzi, W. Lee, "Classification of Packed Executables for Accurate Computer Virus Detection", Elsevier Pattern Recognition Letters, 29(14), pp. 1941-1946, 2008.
18. R. Perdisci, A. Lanzi, W. Lee, "McBoost: Boosting Scalability in Malware Collection and Analysis Using Statistical Classification of Executables", Annual Computer Security Applications Conference (ACSAC), pp. 301-310, IEEE Press, USA, 2008.
19. Protection ID - the ultimate Protection Scanner, available at <http://pid.gamecopyworld.com/>.
20. M. Pietrek, "An In-Depth Look into the Win32 Portable Executable File Format, Part 2", MSDN Magazine, March, 2002.
21. Project Malfease, available at <http://malfease.oarci.net/>.
22. M.G. Schultz, E. Eskin, E. Zadok, S.J. Stolfo, "Data mining methods for detection of new malicious executables", IEEE Symposium on Security and Privacy (S&P), pp. 38-49, USA, 2001.
23. M.Z. Shafiq, S.M. Tabish, F. Mirza, M. Farooq, "A Framework for Efficient Mining of Structural Information to Detect Zero-Day Malicious Portable Executables", Technical Report, TR-nexGINRC-2009-21, January, 2009, available at <http://www.nexginrc.org/papers/tr21-zubair.pdf>
24. M.Z. Shafiq, S.M. Tabish, M. Farooq, "PE-Probe: Leveraging Packer Detection and Structural Information to Detect Malicious Portable Executables", Virus Bulletin Conference (VB), Switzerland, 2009.
25. Symantec Internet Security Threat Reports I-XI (Jan 2002-Jan 2008).
26. F. Veldman, "Heuristic Anti-Virus Technology", International Virus Bulletin Conference, pp. 67-76, USA, 1993.
27. VX Heavens Virus Collection, VX Heavens website, available at <http://vx.netlux.org>.
28. S.D. Walter, "The partial area under the summary ROC curve", Statistics in Medicine, 24(13), pp. 2025-2040, 2005.
29. I.H. Witten, E. Frank, "Data mining: Practical machine learning tools and techniques", Morgan Kaufmann, 2nd edition, USA, 2005.