

Design of Binary File Features for Malware Classification

Jakub Pružinec



Abstract

Rapid and widespread adoption of information technologies lead to userbase diversification and their use among laymen. Due to this, we witness malicious software evolve and grow larger day by day endangering data of billions of users. Today's anti-malware companies seek automated solutions for malware detection. One of possible approaches to malware identification is to use artificial intelligence to classify it. Precision of malware classification is heavily dependent on available information about classified samples - features. Poor design of features may result in wrongly classifying legitimate software as malware, or even worse, to let malware slip by undetected. This article focuses on design, extraction, reliability and efficiency testing of static binary malware features. Moreover, malware feature extraction tool, FileInfo, is innovated. Work is done in cooperation with Avast company, where FileInfo is used in malware clustering system, binary file decompiler and as a general purpose static analysis tool on a daily basis.

Keywords: Binary malware analysis — Static analysis — Classification — Features — Reverse engineering

Supplementary Material: [FileInfo source code](#)

*j.pruzinec@gmail.com, xpruzi02@fit.vutbr.cz, Faculty of Information Technology, Brno University of Technology

1. Introduction

2 Due to modern world digitalization, property, privacy,
3 and identity of people became dependent on informa-
4 tion technologies. Unfortunately, crime has adopted
5 digital nature as well. AV-tests institute recorded over
6 137 millions of malicious software samples in 2018
7 [4]. Cybersecurity became crucial to information tech-
8 nology vendors as they are responsible for safety of
9 their clients. As a consequence of such frequent oc-
10 currence of malicious software, manual analysis of
11 all samples is practically impossible. Therefore, auto-
12 mated malware classification is an absolute necessity
13 when fighting today's malware.

14 Unsurprisingly, it is in great interest of malware

to evade analysis and deceive classifiers to remain un- 15
detected. Evolution of malware and innovation of its 16
obscure evasion methods may result in its missclas- 17
sification. On the other hand, the fact that behavior 18
and shape of malware often differ significantly from le- 19
gitimate software is crucial for malware classification. 20
To classify malware, numerous features are extracted 21
from analyzed samples first. Designed features are 22
general enough to reflect similarity between malware 23
samples and fitting enough to distinguish between ma- 24
licious from legitimate software. What's more, feature 25
design has significant impact on time efficiency of clas- 26
sification algorithms, as it is determined by amount of 27
features taken into consideration (dimensions). Qual- 28

ity of designed features used in a classifier can be estimated by ratio of correctly to incorrectly classified samples.

Despite significant work done in the field of static analysis of malware, most freely available tools are not suitable for malware classification. As a master thesis, Katja Hahn developed a complex malware static analysis tool, *PortEx* [5]. PortEx is a robust open-source library for dissection of binary malware. Unfortunately, PortEx is implemented in Java, thus potentially too slow to be used for malware classification. Further, PortEx supports only Windows binary file format and lacks support of most features presented later in this article.

FileInfo is a complex open-source feature extraction tool developed by Avast company [6]. Contrary to PortEx, FileInfo is implemented in C++ and supports numerous file formats like PE (Windows), ELF (Linux) and MACHO (Apple). Additionally, FileInfo supports extraction of framework based features, e.g. .NET features. Despite its complexity, FileInfo lacks support of features necessary for classification of ever-evolving malware presented hereafter.

This article primarily focuses on (but is not limited to) Windows binary file format static features. Goal of this work is to:

- produce cryptographic and perceptual hashes of program icons
- reconstruct TypeRef tables of .NET binaries and produce their cryptographic hashes
- parse, reconstruct and produce cryptographic hashes of several VisualBasic metadata structures
- compute section and overlay entropy as an indication of packed data
- extract information about product version, supported languages, trademark, copyright, original name and more
- determine addresses of thread-local variable initialization routines

Designed features play a major role when classifying stealthy malware. Malware implementing following analysis evasion methods can be now classified:

- conventional way of importing of external symbols substituted with .NET or VisualBasic symbol importing system
- distortion of icon to corrupt their cryptographic hashes
- packing of sections and overlay to evade static analysis

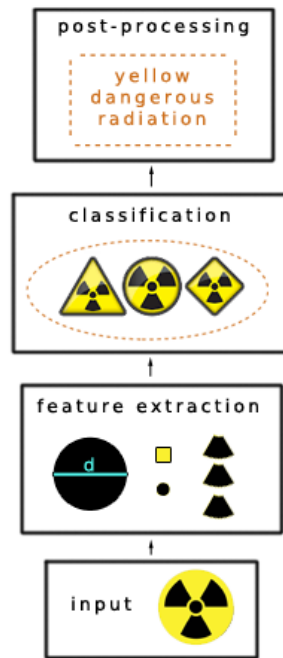


Figure 1. Classification process

Section 2 describes classification process in general, Sections 3, 4 and 5 discuss topic-based features. Additional features are presented in Section 6.

2. Malware Classification

Malware classification process starts with feature extraction from analyzed samples. These features can be of different nature, for example imported symbols or sequence of kernel calls. Afterwards, the classifier assigns classes to input samples based on previously extracted features. As a final step, information about the classes themselves is extracted. Classification process is depicted in Figure 1. Obviously, precision of classification is heavily dependent on extracted features. Features need to be discriminative, meaning the classifier has to be able to distinguish between classes based on them. Poorly designed features may result in classification of legitimate software as malware or the other way around.

Besides proper classification, well designed features may significantly increase time efficiency of classification algorithms. Time complexity of classification algorithms is dependent on number of features taken into consideration.

Features need to be general enough to reflect similarity of malware but should not be too general as it could affect their discriminability. In this context, information about imported libraries is much more valuable than observation that an analyzed sample contains some code. On the other hand, features need to be fitting enough to provide detailed information about



Figure 2. Icon modified with noise to evade classification based on exact match of cryptographic icon hash

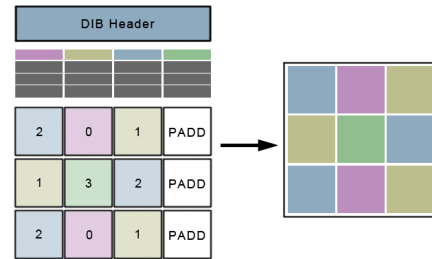


Figure 3. DIB with 4bit color depth parsed into uniform representation

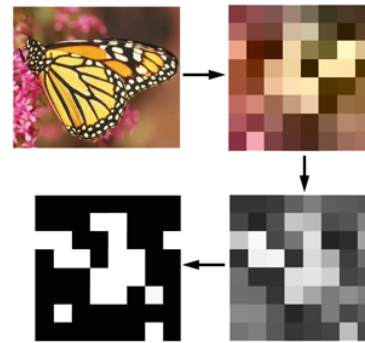


Figure 4. Average hash computation

Figure 4. Image is first resized to 8×8 dimensions to filter out image details. Image is then converted to greyscale and to black and white afterwards. Each of 8 rows of the newly created image is a sequence 8 black (0) or white (1) pixels. Therefore, 8 rows of new image form an 8 byte long average hash. Similarity of two images can be determined by *Hamming distance* of their Average hashes. If the distance is less than 3, images are considered to be similar.

Both cryptographic and perceptual hashes are presented in Listing 1.

Listing 1. Icon hashes

```
resourceTable: {
  "iconAvgHash" : "b7478387b4ffaeff",
  "iconCrc32" : "c6009c34",
  "iconMd5" : "8b6fdb44e0b3e55bf9bc8f-fda1800b79",
  "iconSha256" :
    "4d8b8a948b29bf38cd8f186e75b58-ed5017ed11aac7eb453752181b58f3bea",
  ...
}
```

4. .NET Features

.NET is an open-source framework developed by Microsoft. *.NET* programs are compiled in PE file format as data, metadata and bytecode interpreted by virtual machine.

FileInfo already supports reconstruction of some *.NET* structures but omits information about imported classes. This information is crucial when classifying malware based on imported symbols, because

classified malware. Overfitting features may suppress detection of similarity between analyzed samples. As an example, information about original file name can be made use of, but knowing that third letter of the name is 'X' alone, is useless. Following sections are dedicated to designed features.

3. Icon Features

From 50,000 PE malware samples, over 24,000 contained an icon. Malware authors often add icons to attract and deceive their victims. They often keep icons almost unchanged when updating or modifying malware. Because of this, icon hashes can serve as good features. Cryptographic hashes work well when testing icons for exact match, however a slightest change in icon data results in complete change of its cryptographic hash.

Malware authors are aware of this fact and add noise to their icons as can be seen in Figure 2. To bypass this inconvenience, icons have to be parsed into internal representation. Parsing of icons is a lengthy and rigorous process that will be described in detail in my bachelor thesis. Simply put, in PE file format icons are embedded into a structure called *resource tree*. Despite the fact that PE binaries can contain multiple icons, only one *main icon* is shown in a desktop environment.

Main icon cannot be determined with certainty since the icon to be shown is dependent on desktop environment properties, such as DPI. Main icon detection algorithm will be described in my bachelor thesis as well.

Icon data itself has to be parsed after the main icon has been extracted from resource tree. Icons are stored in *DIB* file format. *DIB* icons of all supported color depths are parsed into uniform representation. Example in Figure 3 demonstrates how a 4 bpp *DIB* icon is parsed into two dimensional pixel array.

When main icon is parsed, one can produce its *perceptual hashes*. Perceptual hashes are hashes representing images in a way that can be tested for similarity. Contrary to cryptographic hashes, slight change of hashed image results only in slight change of its perceptual hash. One such hash is called *Average hash*. Principle of Average hash computation is shown in

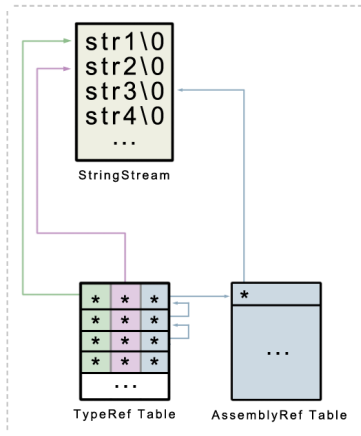


Figure 5. TypeRef table structure

183 .NET does not import most of its external functionality
 184 through conventional *PE import table*. Classes are
 185 imported through so called *TypeRef table*. Each entry
 186 of a TypeRef table may contain information about origin
 187 of an imported class such as *name*, *module type*,
 188 *library name* and *namespace*. Further, .NET classes
 189 may be nested, meaning one class can be defined inside
 190 of another class. On binary level, TypeRef entry of
 191 child class does not reference an external module
 192 but rather TypeRef entry of a parent class. Such
 193 situation is demonstrated on an example in Figure 5.
 194 Reconstruction of a TypeRef table takes several steps:
 195 parsing, linking, and presentation. During the linking
 196 process of a parsed parent/child TypeRef records,
 197 entries are checked for cyclic references. Cyclic
 198 references are exclusively a product of manual
 199 modification of a TypeRef table and may result in
 200 infinite table processing. To prevent this, one of
 201 the references in cycle is simply ignored and remains
 202 detached. After reconstruction of a TypeRef table,
 203 its cryptographic hashes are produced.

204 Listing 2 shows the first element in a reconstructed
 205 form of a TypeRef table.

Listing 2. TypeRef table features

```

206 "typeRefTable" : [
207   { "libraryName" :
208     "System.Runtime",
209     "name" :
210       "CompilationRelaxationsAttribute",
211     "namespace" :
212       "System.Runtime.CompilerServices"},
213 ]

```

214 5. VisualBasic Features

215 Despite the claim that *VisualBasic* programs “benefit
 216 from security” [7], today they are almost exclusively
 217 created by malware authors. Same as .NET, VisualBasic
 218 is interpreted by a virtual machine, thus implements

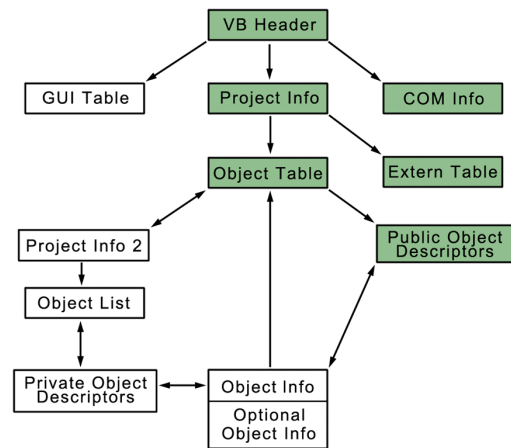


Figure 6. VisualBasic

its own importing system besides *PE import table*. 219

VisualBasic format is closed-source, therefore it 220
 is easier for VisualBasic malware to hide its inten- 221
 tions. On the other hand, VisualBasic applications 222
 are rich in metadata. Alex Ionescu has done remark- 223
 able work on reversing the metadata [8] and therefore 224
 lots of previously hidden VisualBasic malware can be 225
 now classified. Metadata includes information about 226
project names, *identifiers*, *languages*, *imported func-* 227
tions, *objects* and more. This information is spread 228
 across numerous structures referencing each other as 229
 shown in Figure 6. 230

To extract the necessary data, a top-down parser 231
 is implemented. For the sake of receiving names of 232
 imported functions, an external table is reconstructed. 233
 Further, an object table is reconstructed to obtain names 234
 and method names of implemented objects. External 235
 table and object table are hashed and all auxiliary fea- 236
 tures are presented. 237

Small excerpt of VisualBasic features is presented 238
 in Listing 3. It can be assumed that malware targets 239
 Facebook users and is written by an author of Swedish 240
 origin. 241

Listing 3. VisualBasic features excerpt

```

// imported functions 242
"externs" : [ 243
  { "apiName" : "PlgBlt", 244
    "moduleName" : "gdi32" }, 245
], 246
// implemented objects 247
"objects" : [ 248
  { "name" : "frmMain" 249
    "methods" : [ 250
      "C_Mutex", 251
      "BROWSER_FB_DocumentComplete", 252
      "BROWSER_FB_OnQuit", 253
      "FACEBOOK_START" }, 254
    ], 255
  ],

```



```

256 // additional
257 "projectPath" :
258   "C:\\Users\\Admin\\Desktop_old\\
259   Blackshades project\\Blackshades
260   NET\\server\\server.vbp",
261 "projectPrimaryLCID" :
262   "English - United States",
263 "projectSecondaryLCID" :
264   "Swedish - Sweden"

```

```

// thread-local initialization routines
"threadLocalInitializers" : [
  "0x401060",
  "0x4010a0"
],

```

Besides the data necessary for proper program execution, compilers often add additional information about the product. Such information can be retrieved from *VersionInfo* resource tree entry in PE file format. This entry holds information about *supported languages, legal copyright, original file name, version, timestamps* and more.

Part of information extracted from *VersionInfo* of a malware is to be found in Listing 6. This information suggest that the author was of German origin. The author probably instrumented a Firefox executable built on 2010/09/14.

6. Other Features

One of commonly used methods to avoid static analysis malware authors use is to *pack* their program. Packed binary contains compressed data and code that are decompressed during runtime by decompression routines. This is particularly inconvenient, because extraction of static features becomes very hard or practically impossible. On the other hand, packed binary can be a strong indication of malicious intentions. For this reason, entropy of sections and overlay is computed. High data entropy indicates compression, thus can serve to detect packed data. Contrary, low entropy indicates small data diversity and can be used to detect blank data sections.

Listing 4 demonstrates detection of section compression status.

Listing 4. Section entropy

```

281 "sections" : [
282   // packed
283   { name: ".text",
284     entropy: "7.8632" },
285   // normal
286   { name: ".rodata",
287     entropy: "4.3242" },
288   // empty
289   { name: ".fini_array",
290     entropy: "0.8632" },
291 ]

```

Thread-local data initialization routines are abused by malware authors to evade static analysis. Thread data is stored in dedicated directory in PE file format. This data needs to be initialized before entry point execution. In other words, thread-local data initializers are called before the *main()* function. Malware often implements its malicious behavior in one of thread-local data initializers for this reason. Intuitive investigation of *main()* function is useless in this case. What's more, many debuggers set breakpoints on entry point, therefore thread-local data initialization routines may infect host machine before the analyst has a chance to intervene.

Addresses are presented as shown in Listing 5.

Listing 5. Thread-local data initializers

Listing 6. VersionInfo features

```

"versionInfo" : {
  "languages" : [
    { "codePage" : "utf-16",
      "lcid" : "German - Germany" }
  ],
  "strings" : [
    { "name" : "CompanyName",
      "value" : "obama" },
    { "name" : "OriginalFilename",
      "value" : "my_st0re_lo-
      ader____.exe" },
    { "name" : "ProductName",
      "value" : "Firefox" },
    { "name" : "BuildID",
      "value" : "20100914121323"
    }
  ]
}

```

7. Conclusions

Design of features and their impact on malware classification have been discussed in this article. Further, some static analysis evasion approaches and methods to deal with them were presented. Following contributions were made:

- design and extraction of .NET TypeRef table features
- design and extraction of icon features
- design and extraction of VisualBasic metadata features
- design and extraction of VersionInfo features
- design and extraction of entropy and thread-local directory features

So far, only icon and .NET features were integrated into Avast malware clustering system. In Table 1 and

357 Table 2 are statistics of clusters classified as malware
 358 based solely on a given feature. Beside the size of a
 359 cluster, detected malware ratio can be seen. Further,
 360 malware detection ratio of 1000 randomly selected
 361 samples from clusters analyzed by ESET and Kaspersky
 antivirus software is present.

| N samples | Detected | ESET | Kaspersky |
|-----------|----------|------|-----------|
| 406 | 97% | 95% | 100% |
| 221 | 96% | 100% | 100% |
| 50 | 96% | 100% | 100% |

362 **Table 1.** Icon MD5 based clusters

| N samples | Detected | ESET | Kaspersky |
|-----------|----------|------|-----------|
| 1.3M | 96% | 100% | 100% |
| 63K | 98% | - | 82% |
| 13K | 98% | 99% | 100% |

Table 2. TypeRef MD5 based clusters

363 In close future, newly designed features imple-
 364 mented in FileInfo will be integrated into the cluster-
 365 ing system and RetDec decompiler developed by Avast
 366 company [6]. FileInfo implements an in depth analysis
 367 of the PE file format prevalently and soon it will be
 368 forced to respond to increasing occurrence of Linux
 369 and macOS malware with extraction of new features
 370 regarding ELF and MACHO file formats.

371 Further, FileInfo has great potential for improve-
 372 ment of data recognition features, such as overlay for-
 373 mat detection. Besides that, features regarding file for-
 374 mat violations should be implemented, as they serve
 375 as good indication of malicious behavior. Some of
 376 the problems mentioned here will be addressed in my
 377 bachelor thesis soon.

378 Acknowledgments

379 I would like to thank my supervisors Doc. Dr. Ing.
 380 Dušan Kolář, Ing. Marek Milkovič and Ing. Jakub
 381 Křoustek, PhD for help and support throughout the
 382 whole project.

- [1] MBIS. Cyber forensic - scientific inves- 384
 tigations. [https://www.mbis-inc.net/](https://www.mbis-inc.net/cyber-forensic-services.html) 385
[cyber-forensic-services.html](https://www.mbis-inc.net/cyber-forensic-services.html). [On- 386
 line; accessed April 6, 2019]. 387
- [2] Michael Sikorski and Andrew Honig. *Practical 388*
malware analysis: the hands-on guide to dissect- 389
ing malicious software. no starch press, 2012. 390
- [3] Warren Perez Araya. Is malware analy- 391
 sis right for your business? [https://](https://securityintelligence.com/is-malware-analysis-right-for-your-business/) 392
[securityintelligence.com/is-](https://securityintelligence.com/is-malware-analysis-right-for-your-business/) 393
[malware-analysis-right-for-your-](https://securityintelligence.com/is-malware-analysis-right-for-your-business/) 394
[business/](https://securityintelligence.com/is-malware-analysis-right-for-your-business/), 2018. [Online; accessed April 6, 395
 2019]. 396
- [4] Av-test institute: Malware statistics 2018. 397
[https://www.av-test.org/en/](https://www.av-test.org/en/statistics/malware/) 398
[statistics/malware/](https://www.av-test.org/en/statistics/malware/). [Online; accessed 399
 April 6, 2019]. 400
- [5] Katja Hahn. Robust static analysis of portable 401
 executable malware. *Mater Thesis, HTWK Leipzig*, 402
 2014. 403
- [6] J. Křoustek, P. Matula, and M. Milkovič. An 404
 open-source machine-code decompiler. [https://](https://retdec.com/static/publications/retdec-slides-recon-2018.pdf) 405
[retdec.com/static/publications/](https://retdec.com/static/publications/retdec-slides-recon-2018.pdf) 406
[retdec-slides-recon-2018.pdf](https://retdec.com/static/publications/retdec-slides-recon-2018.pdf), 2018. 407
 [Online; accessed April 6, 2019]. 408
- [7] Microsoft corporation: Visual basic guide. 409
[https://docs.microsoft.com/en-us/](https://docs.microsoft.com/en-us/dotnet/visual-basic/) 410
[dotnet/visual-basic/](https://docs.microsoft.com/en-us/dotnet/visual-basic/), 2018. [Online; 411
 accessed April 6, 2019]. 412
- [8] Alex Ionescu. *Visual Basic: Image Internal Struc-* 413
ture Format. Relsoft Technologies, 2004. 414