# On the feasibility of binary authorship characterization

Saed Alrabaee [a,*,1], Mourad Debbabi [b], Lingyu Wang [b]

[a] Information Systems & Security, United Arab Emirates University, United Arab Emirates
[b] Computer Security Laboratory, Concordia Institute for Information Systems Engineering, Concordia University, Canada

## ARTICLE INFO

*Article history:*

## ABSTRACT

This work aims to develop an automatic tool that can perform the laborious and error-prone reverse engineering task of binary authorship characterization, i.e., determining clues related to the author(s) of a piece of binary code. Software code written by human programmers reflects the author's educational background, level of expertise, and coding traits. Accordingly, these may be characterized by identifying meaningful features and examining them. Binary authorship characterization reveals information that can be extremely useful for security applications such as digital forensics, malware triage, and binary vulnerability tracking. This paper proposes a system, `BinChar`, that capture various aspects of author style, including code trait characteristics, code structure characteristics, and code behavior characteristics. For the purpose of detection, a Convolutional Neural Network (CNN) is used. The results generated by the CNN are evaluated more precisely using Bayesian calibration. We tested `BinChar` in identifying the characteristics of the authors of program binaries. Also, we applied it to almost 500 GB of malware samples provided by the Kaggle Microsoft Malware Classification Challenge, to demonstrate that `BinChar` is an appropriate tool for characterizing malware families. As an illustration, we report a case study in which we determine the author characteristics of the Mirai botnet and compare them with the author characteristics of 360,000 malware samples.

## Introduction

When analyzing malware binaries, reverse engineers often pay special attention to their characterization for several reasons. First, reports from anti-malware companies indicate that finding the similarities between malware code characteristics can aid in developing profiles for malware families (Techniqal report and Resource, 2012). Second, recently released reports by Citizen team (Big Game Hunting, 2015; Citizen Lab, 2015) show that malware binaries written by authors having the same origin share similar characteristics. Third, many malware packages could have been written only by authors with a special level of expertise and special knowledge for dealing with specific resources; an example is SCADA system malware. This insight provides a critical clue for the extraction of information about the functionality of a malware binary. Fourth, although obfuscation techniques may be applied before the malware is released and may modify the code significantly, it is still desirable to determine which obfuscation

techniques and tools have been used. Last, clustering binary functions based on a common origin may help reverse engineers identify the group of functions that belong to a particular malware family or decompose the binary based on the origin of its functions.

The ability to conduct these analyses at the binary level is especially important for security applications because the source code for malware is not always available. However, in automating binary authorship characterization, two main challenges are typically encountered: the binary code lacks many abstractions (i.e., function prototypes) that are present in the source code; and the time and space complexities of analyzing binary code are greater than those of the corresponding source code. Although significant efforts have been designed to develop automated systems for source code authorship characterization (Caliskan-Islam et al., 2015; Frantzeskou, 2004; Taylor et al., 2008; Woldring et al., 2016), these often depend on features that will likely not be preserved in the strings of bytes representing executable file after the compilation process, such as variable and function naming, original control and data flow structures, comments, and space layout.

To the best of our knowledge, there have been no attempts to characterize the authors of program binaries. Nonetheless, a few approaches to binary authorship attribution have been proposed, but they typically use machine learning algorithms to extract unique patterns for each author and then compare a given target

---

* Corresponding author.
  *E-mail address:* salrabaee@uaeu.ac.ae (S. Alrabaee).
[1] Majority of the work was conducted during Saed's work at University of New Haven.

binary against such patterns to identify the author (Alrabaee et al., 2014; Alrabaee et al., 2018; Caliskan-Islam et al.; Meng, 2016; Rosenblum et al., 2011).

These approaches cannot be applied directly to binary authorship characterization because of the following limitations: the chosen features are generally not related to author style but rather to functionality; they are not applicable to real malware; and dealing with the binary authorship characterization problem requires that the chosen features have the power to detect, for example, author styles, the compilers used, the free packages reused, the functionality, the implementation frameworks, and the binary timestamps. More recently, the feasibility of authorship attribution for malware binaries was discussed at the BlackHat conference (Big Game Hunting, 2015). It was concluded that a set of features can be employed to group malware binaries according to authorship characterizations. However, the process is not automated and requires considerable human intervention.

**System overview.** To address the aforementioned limitations, this paper presents an innovative system, `BinChar`, that describes the characteristics of programmers according to their educational background, level of expertise, and coding traits. To achieve this, we have defined a new set of features that extracts authorship attribution characteristics. These features are extracted from different levels: the level of the basic blocks, the level of the function bodies, the program level, and the file level. Based on them, our system is able to detect structural, optimization, knowledge, expertise, and code trait characteristics, and also overall characteristics that is resulted from the aforementioned characteristics together. All the extracted features are passed to the CNN. We observe note some attractive characteristics of neural networks such as they can learn end-to-end, where each stage is trained simultaneously to achieve the end goal (Shin et al., 2015). Also, the nodes of CNN can act as filters over the input space and can discover the strong correlation in the binary code (lab, 2014). Third, CNN is considered a fast neural network in classification process when it is compared with other neural networks such as recurrent neural network (Wei et al., 2016). The results obtained from CNN is passed to Bayesian calibration to precisely check the correctness of CNN.

**Contributions.** Our contributions are summarized below.

- We designed a new set of features that make *BinChar* accurate and efficient, enabling it to characterize the authors of program binaries with high speed while tolerating the noise injected by code transformations arising from the use of different compilers and optimization speed levels. The experimental results show that our system is able to cluster the samples according to similarities in authorship characteristics with a precision of over 95%.
- We investigated the effectiveness and the power of CNNs in the context of binary authorship characterization. To the best of our knowledge, this is the first work in which CNNs are used for binary authorship characterization. Further improvement is achieved by performing Bayesian calibration, which reduces the rate of false positives to 0.02%.
- We used `BinChar` to extract author characteristics from Mirai botnet binaries and compared the results with those for 360,000 malware samples collected from various sources. Finally, we report the authorship characteristics common to the Mirai botnet and other families.

## Preliminaries

### Threat model

`BinChar` is designed to assist, instead of replacing, reverse

engineers in various use cases of characterizing the author of program binaries, such as forensic analysis (e.g., linking a new malware to previously known malware or malware author(s), clustering a group of malware based on common characteristics, and finding co-authorship for the malware binary code (e.g., to determine "influencers" in the code community), and software copyright infringement analysis (e.g., detecting borrowed code fragments inside a program binary based on authorship characteristics). Therefore, the focus of `BinChar` is not on general reverse engineering tasks, such as unpacking and de-obfuscating malware samples (although `BinChar` leverages some existing tools to de-obfuscate and unpacked such binaries), but rather on detecting and determine authorship characteristics clues such as using advanced resources, particular compiler, os, specific code traits, memory allocation habits, c&c commands, etc. We just pay attention to the author characteristics which can be common among different set of authors). Moreover, the authors of a specific malware family, they have to share common knowledge and expertise to deal with a particular environment (Nagano and Uda, 2017; Wagner et al., 2017). In designing the features and methodology of `BinChar`, we have taken into consideration some potential countermeasures. In particular, `BinChar` assumes the adversary may attempt to evade detection through the following.

- The adversaries may apply refactoring techniques, e.g., when a malware author aims to defeat forensic analysis by modifying his/her own code, or when an adversary attempts to modify borrowed code written by other authors in order to evade copyright infringement detection.
- The adversaries may apply obfuscation techniques on binary files to alter its syntax, e.g., when a malware author wants to defeat anti-virus signature-based detection.
- Since a program can be significantly altered by simply changing the compilers or their settings, the adversary may make such changes to evade detection.

We will show in later sections how `BinChar` may survive the aforementioned threats. Simply, the features of `BinChar` have been designed to determine binary authorship characteristics at multiple abstraction levels, which makes it harder for adversaries to evade system detection (Big Game Hunting, 2015). In addition, an operational solution is to customize and enrich the list of features employed by `BinChar` based on the actual use case and learning data, which will not only make it much more difficult for adversaries to hide all of their binary code characteristics such as code traits, but will also improve accuracy.

### Usage scenarios

The interest of the `BinChar` framework is to determine the authorship characteristics of an anonymous piece of binary code. Additionally, given code that is written by multiple authors, it is required to determine which characteristics (e.g., coding habits) of the code are attributed to which author. It is assumed that a set of binary code samples is available, where the code is either labeled with an author(s) from a set of known candidate authors, or from other unknown authors. Given an anonymous piece of binary code, `BinChar` converts the code into a set of features that are consequently used to attribute the author(s) according to their characteristics. `BinChar` also labels the compiler-related functions in addition to discovering the reused functions and third-party libraries. This is useful for a number of applications, as listed in what follows.

**Software Infringement**. A set of candidate authors is clustered based on previously collected malware samples, online code

repositories, etc. There are no guarantees that the anonymous author is one of the candidates, as the test sample may not belong to any known authors. However, it is not our intention to identify the author, but rather to characterize the author. Finally, it may be suspected that a piece of code is not written by the claimed author, but yet there are no leads as to who the actual author may be (Caliskan-Islam et al., 2015). For this reason, we may compare the programs written by the claimed author and measure the degree of similarity in terms of binary code characteristics.

**Forensic Investigation**. FireEye (Moran and Bennett, 2013) discovered that malware binaries share the same digital infrastructure and code (for instance, the use of certificates, executable resources and development tools). FireEye investigators eventually noticed that malware binaries of the same previously-discovered infrastructures are written by the same group of authors. In such cases, training on such binaries and some random authors' code may offer vital assistance to forensic investigators. In addition, testing recent pieces of malware binary code using confidence metrics would verify if a specific author is the actual author.

*System overview*

The architecture of `BinChar` is illustrated in Fig. 1. As shown, the four main components are: (i) Preprocessing component where PEfile (file, 2012) is employed to check if the binary file is packed. If it is packed, the corresponding unpacker, such as UPX, is used to unpack the binary file and pass the unpacked binary file to the disassembler tools such as IDA Pro (HexRays, 2011). (ii) Feature extraction which deals with a different set of features that are related to the authorship characterization. This component is able to accept either an executable file or an assembly file as input. In the case of executable file, BinChar converts executable file to a gray-scale image. Then numerous features are extracted from this image. After that we rank these features by employing mutual information. These top ranked features are passed to detection component. Also, BinChar uses executable file as input to LLVM in order to optimize the code and compare it with the original executable. In the case of assembly file, BinChar extracts the following static features: (1) merging Annotated Control Flow Graph (ACFG) with Data Flow Graph and the new graph is decomposed then into graph truss. (2) We extract the longest path and then apply statistically analysis it. (4) Strings are extracted that can disclose the binary provenance such as the source compiler. The extracted featured are hashed by LSH algorithm. (iii) The fixed size of hashes resulted by LSH are passed to the third component. The CNN is trained by set of fixed size of hashes. Once we have a target binary, CNN will return a set of candidates according to the training hashes. These candidates are passed to Bayesian model calibration to seek to provide users with accurate probabilities that a given binary file. This component is used to identify the author of program binary if the characteristics are the same. (iv) The final component represents the outputs of `BinChar`. `BinChar` clusters function or programs together according to the similarity of the author's characteristics of the binary code. Standard clustering algorithm (k-means) (Farnstrom et al., 2000) is used for this purpose.

As shown in Fig. 1, the author repository is connected to this component since our features can reveal a significant information such as the compilers, resources used, level of expertise, education background, API calls, and some code traits like the way of configuration.

`BinChar`: **design overview**

In this section, we describe our system in detail.

*Feature engineering*

*Optimization characteristics*

We designed a set of features based on existing tools (Tate et al., 2009; Tristan et al., 2011; Schkufza, 2015; Sharma, 2016). We leveraged these tools to construct a data flow graph incorporating data and control dependencies. The graph reveals the inline functions, making it easy to count them. Then we constructed a semantic graph and normalized it by applying the rules introduced in (Tristan et al., 2011); these are related to code optimization. The first rule is constant folding. For instance, the 3-node subgraph representing the expression `1 + 4` is rewritten into a single node to represent the value `5`. The second rule concerns side effects. The authors of (Tristan et al., 2011) employ state variables to detect the dependencies between instructions. They define a state variable for memory states: every instruction contains of a memory register and its usual parameters. This additional register forces a dependency between, for example, a load instruction and the preceding store instructions (Tristan et al., 2011), retaining all the relevant information about each register. The third rule concerns aliasing information that relaxes the strict ordering of instructions imposed by the transformation. For instance, the pointers in LLVM returned by `alloc` never alias with one another. Hence, if two load instructions from a memory involve a store to the same pointer, then the load can be simplified to a single instruction. The fourth rule deals with loops. The authors came up with a method for placing gates within looping control flow, including breaks, continues, and returns from within a loop. Once all the rules were applied, we conducted statistical analyses to compare the original code with the optimized code, including the mean, variance, standard deviation, and mean square error.

*Structural characteristics*

This subsection will show the merging of the Annotated Control Flow Graph (ACFG) with the Data Flow Graph (DFG) into a single graph, the "author style" graph. The reasoning underpinning this merger is that the ACFG may offer insight into the implementation of the task by the author, and the DFG will show the ways in which the author has manipulated variables and will also show how they have employed small functions, which the compiler generally inlines while optimizing.

*Author Style Graph.* We constructed a semantic control flow graph named the *Annotated Control Flow Graph* (`ACFG`). This represents an abstracted form of the CFG, with multiple criteria informing the
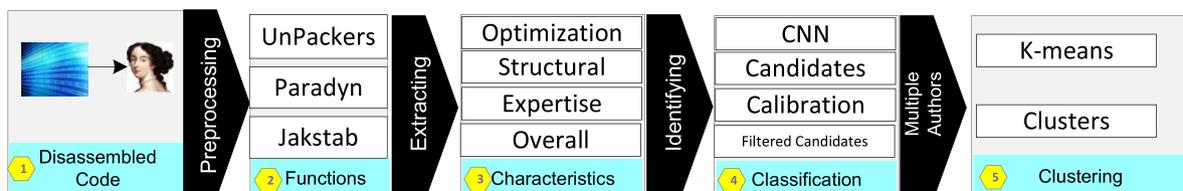


**Fig. 1.** *BinChar* architecture.

generalization of specific types of features from the CFG. There follows a summary of the building process for the ACFG. The system has an input of a CFG and calculates how frequently opcodes occur in the instruction groups. It then assigns assembly instructions to one of six categories (Rahimian et al., 2015), illustrated in Table 1.

The resulting ACFG describes structural characteristics. The next step is to complement it with the Data Flow Graph (DFG). Then, data flow dependencies can be incorporated to make it possible to use coarse reasoning about the program control flow and data dependencies to infer how the author manipulates program variables and how variable relations are built. Let $R_k/W_k$ denote registers or memory that are read or written by instruction $I_k$. If $i_1$ and $i_2$ are instructions belonging to $I$, then we define the following possible dependencies: $i_1$ writes something which will be read by $i_2$; $i_1$ reads something before $i_2$ overwrites it; and $i_1$ and $i_2$ both write the same variable. Recognizing that the ACFG and DFG provide complementary views about the implementation of a specific task by emphasizing different aspects of the underlying authorship characteristics, we combine them into a joint data structure to facilitate more efficient graph matching between different binary codes to help identify similarities in authorship characteristics. The result is the Author Style Graph (ASG).

**Definition 1.** The Author Style Graph $G = (N,V,\zeta,\gamma,\vartheta,\lambda,\omega)$ is a directed attributed graph, where $N$ is a set of nodes, $V \subseteq (N \times N)$ is a set of edges and $\zeta$ is an edge labeling function which assigns a label to each edge: $\zeta \to \gamma$, where $\gamma$ is a set of labels. Finally, $\vartheta$ is a data dependency function which labels each node $n \; \varepsilon \; N$ based on the data dependency rules function $\lambda$.

### Expertise characteristics

Typically the path in the CFG has been shown as a robust feature (Huang et al., 2017). Further, to find matched paths can be considered as an alignment problem where dynamic programming can be applied (Huang et al., 2017). Longest path reflects the author traits in implementing tasks or using nested loops. Thus, we choose the longest path. We use depth first search to traverse the CFG, and then choose the path with the largest number of nodes. Once the longest path is constructed, we extract two categories of features: statistical features (cyclomatic complexity) and dynamic features (execution traces).

Cyclomatic complexity (McCabe, 1976) computes the quantity of linear independent paths to represent the complexity of a code in the form of a graph metric. The complexity C of the longest path of CFG is as follows: $C = E + N + 2P$, where $E$ is the edge quantity, $N$ is node quantity, and $P$ is the quantity of connected components in the longest path.

### Overall characteristics

The binary file is converted into a set of bytes, and the set is then transferred to a matrix. This is a two-dimensional matrix with the width fixed at $d$; to be more precise, the first $d$ bytes are placed on the first row of the matrix, with the $nth$ group of $d$ bytes moving to the $nth$ matrix row. This methodology has similarities with the

**Table 1**
Patterns for annotation (Rahimian et al., 2015).

| Feature | Description | Example |
|---|---|---|
| DTR & STK | Data Transfer and Stack | push, mov, etc |
| ATH & LGC | Arithmetic and Logical | add, xor, etc |
| CAL & TST | Call and Test | call, cmp, etc |
| REG & MEM | Register and Memory | esi, [esi+4] |
| REG & CONST | Register and Constant | esi, 30 |
| MSC | Milestones | MEM and Const |

techniques for visualizing malware outlined in (Kirat et al., 2013; Nataraj et al., 2011, 2013).

### Feature hashing

As our tool employs a number of features that have been extracted using different formats (instruction traces, graphs, numeric, etc), we chose LSH so that features could be unified prior to them being fed into CNN. *minhash* (with $K$ unique hash functions) are employed as a signature set for the introduction of used features, and the minhash values will be divided by LSH to form a signature matrix of $I$ bands, each containing $r$ rows (Karbab et al., 2016). For calculation of the quantity of minhash values for a particular band, the number of minhashes is divided by the number of bands ($K/I$). The number of rows and the number of featured minhash signatures will be equal.

## System detection

### Convolutional Neural Network

Convolutional Neural Networks (CNNs) are one type of Neural Networks that have shown efficiency in various areas such as image processing. Generally, each input node is connected to each output node in the next year. In CNN, it is not the case where convolution is used over the input layer to compute the output. This process leads in local connectivity, each input region is connected to a node in the output. Many filters are applied at each layer and combine their results. After convolution step, the subsequent step is called ReLU which stands for Rectified Linear Unit and is a non-linear operation. Its output is given by: Output = Max(zero, Input). This step is used to replace any unwanted features by zero. The purpose of ReLU is to introduce non-linearity in CNN. This operation is called activation function. Besides, there are other layers are called pooling layers. These layers are inserted periodically between convolution layers. The main goal of pooling layers is to reduce the spatial size of the representation, to minimize the number of parameters and computation in the network, and also to dominate overfitting. Finally, the last layer is then a classifier that uses these high-level features. It has been show that CNN can be more accurate and efficient to train if they contain shorter connections between layers close to the input and those close to the output (Huang et al.). Consequently, we have adopted their CNN algorithm for binary authorship characterization. Their algorithm follows a feed-forward fashion by connecting each layer to every other layer. Where the inputs used for each layer, are the feature-maps of all preceding layers. Besides, its own feature-maps are used as inputs into all subsequent layers (Huang et al.). This has one main advantage by alleviating the vanishing-gradient problem.

Let us examine a feature vector $v_0$, which we transfer via a convolutional network. This network comprises $L$ layers, all of them implementing nonlinear transformations $T_l$ in which the layer is indexed by $I$. $H_l$ may be rectified linear units(ReLU) (Glorot et al., 2011; Maas et al., 2013). The output of the $i^{th}$ layer is denoted as $v_i$.

**Identity function**. For traditional CNN forward networks, the $l^i$ output layer as input to $(l+1)^i$, which creates this layer transition: $v_i = H_l(v_{i-1})$. A skip connection with a transition layer is added in order to bypass the non-linear transformation and create an identity function: $v_i = H_l(v_{i-1}) + v_{i-1}$. This process has the advantage that there is a direct flow gradient between later layers and earlier layers via the identity function.

**Dense connectivity**. Improvements should be made regarding the way information flows between layers to make it more accurate and efficient. We therefore have employed the algorithm proposed by (Huang et al.). This introduces a direct connection between any

layer and every subsequent layer. Thus, the $i^{th}$ layer is in receipt of the feature-maps for every preceding layer, $v_0, \cdots, v_{i-1}$, as input: $v_i = H_i([v_0, v_1, \cdots, v_{i-1}])$, where $[v_0, v_1, \cdots, v_{i-1}]$ refers to the merging of the feature-maps created in layer $0, \cdots, i-1$.

**Activation function**. To achieve this purpose we have used a rectified linear unit (ReLU) (Maas et al., 2013). Separate computations are undertaken for each layer within this unit. If this unit was absent, the layer would be an affine function.

**Pooling layer**. This means that pooling the layers that are responsible for changing the size of feature-maps is an important element of convolutional networks. To allow for pooling in the network, it is divided into different dense components following the architecture in (Huang et al.). The layers and their connectivity are illustrated in Fig. 2. As Fig. 2 shows, transition layers exist between blocks that are responsible for pooling and convolution. Our model employs the same settings employed by (Huang et al.). A batch normalization layer is employed with a 1x1 convolutional layer and then a 2x2 average pooling layer (Huang et al.).

*Calibration model*

Motivated by (Kolosnjaji et al., 2016), we are interested in using their Bayesian calibration model. The main goal of this model is to deliver users with accurate probabilities that a given file is how similar to another file. This calibration model combines our firs probability threshold value (represented as the ratio of author $X$ to all other authors in our dataset) and the other information about our CNN's error profile against test data. In what follows we briefly introduce our adapted model for adjusting the CNN similarity score to reflect the true author label score, given this qualitatively assumed ratio of author to all other authors.

Considering $0 \le c \le 1$ be some score given by the CNN, which indicates the degree to which a CNN measures a binary how similar is to other author of program binaries, with 0 being completely not related to any author, and 1 being certainly related to that author. Our goal is to convert this number into a "calibration" score, which will provide the user a measure of how likely that the author of target binary is accurately related to the authors of candidates. To achieve this, we define the calibration score as the probability that the author of program binary will accurately related to the resulted candidates by CNN, $P(R = x_i | C = c)$, given the score $c$, and list of candidates $R = x_1, x_2, \cdots x_i, x_{i+1}, \cdots x_k$, where k is the threshold value that represents the number of candidates. By using Bayes' law we have:

$$P(x_i|c) = \frac{p(c|x_i)P(x_i)}{p(c)} =$$

$$\frac{p(c|x_i)P(x_i)}{p(c|x_i)P(x_i) + p(c|u)P(u)}$$

where $p$ is the pdfs, we suppose that pdfs for the list of candidates and target scores for CNN, $p(C = c | R = u)$, where u is the set of all candidates except the target file, which means $u = \{x_1, x_2, \cdots x_i, x_{i+1},$ $\cdots x_k\} - x_i$.

In order for the probabilities to total 1, a constraint is employed that provides a final value for a recalibration score, related to PDFs and the probability of the target author being in the list of candidates. Our problem is defined through the CNN's predicted PDF for target binary and candidate list $u$. For deriving the PDFs, we employed a non-parametric approach, such as the kernel density estimator (KDE) (Saxe and Berlin, 2015), with a value of PDF being approximated given C through a weighted neighborhood average. As the pdfs may only assume values in [0,1], mirroring of the samples to the left of 0 and the right of 1 takes place, prior to the computation of the estimated pdf value for a specific point.

*Clustering similar functions*

For the clustering process, we use a standard clustering algorithm (k-means) (Farnstrom et al., 2000) to group functions with similar author characteristics attributes $(v_{n_1}, \ldots, v_{n_z})$ into $k$ clusters $S = S_1, \ldots, S_k$ and $(k \le |F_{P_1}| + |F_{P_2}|)$ to minimize the intra-cluster sum of squares. For more details, we refer the reader to (Farnstrom et al., 2000).

## Evaluation

In this section, we present the evaluation results for the possible use cases described earlier in this paper. Section 5.1 shows the setup of our experiments and provides an overview of the data we collected. The main results on authorship attribution characterization and identification are then presented. Also, we have studied the impact of different CNN parameters on the *BinChar* accuracy. Finally, *BinChar* is applied to real malware binaries and the results are discussed.

*Implementation environment*

The described binary feature extractions are implemented using separate python scripts for modularity purposes, which altogether form our analytical system. For CNN setup, we first use a convolution with 16 output channels is performed on the input feature vectors before the first dense block. We use kernal size $3 \times 3$ for convolutional layers. We follow zero-padded for inputs to keep the feature-map size fixed (Huang et al.). We use $3 \times 3$ convolution followed by $4 \times 4$ average pooling as transition layers between two contiguous dense blocks. Further, the global average pooling is excuted at the end of the last dense block. Then a softmax classifier is attached. The feature-map sizes in the two dense blocks are $128 \times 128$, and $64 \times 64$, respectively.

*Dataset*

The used dataset is consisted of several files from different sources, as described below: i) GitHub (The GitHub repository, 2016); ii) Google Code Jam (The Google Code Jam, 2008–2015); and iii) a set of known malware files representing a mixture of 1500
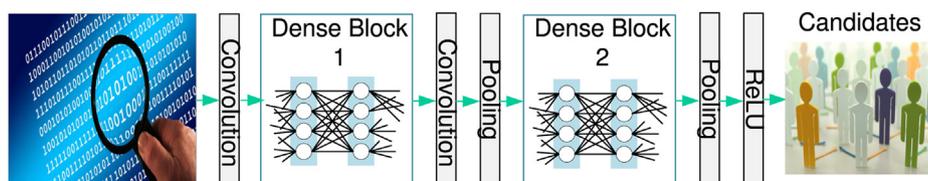


**Fig. 2.** Our CNN architecture with deep DenseNet (Huang et al.) with two dense blocks.

different families including the nine families provided in Microsoft Malware Classification Challenge (Big Game Hunting, 2015). Statistics about the dataset are provided in Table 2.

*Dataset compilation*

The ultimate dataset is compiled with different compilers and compilation settings to measure the effects of such variations. We use g++, Clang, GNU Compiler Collection's gcc, ICC, as well as Microsoft Visual Studio (VS) 2010, with different optimization levels.

*Accuracy*

The purpose of this experiment is to evaluate the accuracy of characterizing the author of in program binaries.

**Evaluation Settings.** The method used in this evaluation employs random subset testing as well as standard ten-fold cross-validation, on the basis of the experiment: To classify the whole data set, we have employed ten-fold cross-validation. We undertook evaluation of the *BinChar* system by employing the datasets referred to in Section 5.2. The data is separated at random into ten sets, with one set being kept as the testing set, with the other nine being training sets. *BinChar* For comparison with extant methodologies, we apply precision recall $R$ and $P$ measures. As the target application domain for *BinChar* is far less sensitive to false negatives than to false positives, the following F-measure is employed:

$$F_1 = 2 \ \frac{P \cdot R}{P + R} \tag{1}$$

**BinChar Accuracy.** We first investigate the accuracy of our proposed system in identifying the author of program binaries based on author characteristics. The results are reported in Table 3. The highest accuracy obtained by our tool is 0.94 when all characteristics components are together. Further, we can observe that the expertise characteristics return the highes accuracy of 0.93. This is due to the fact that the author may use his expertise to implement a specific task. For instance, the author may use specific package or advance resources to reduce the execution time. As mentioned earlier, we use static and dynamic feature to detect the author expertise characteristics. However, the optimization characteristics return an accuracy of 0.81. Here, we optimize the original code and then we compare them in terms of statistical analysis. We observe through our experiments if the author tries to optimize the written code, then these characteristics may not help fully to identify the characteristics of the author.

*False positive rate*

We investigate the false positives in order to understand the situations where *BinChar* is likely to make incorrect attribution decisions. The average of false positive rate is 0.02%. It is very low and could be neglected. The reason behind the low false positive rate is that BinChar uses stratified detections system. We have observed that the false positives rate for google dataset is the

**Table 2**
Statistics about the binaries used in the evaluation.

| Source | # of authors | # of programs | # of functions |
|---|---|---|---|
| GitHub | 3550 | 12,910 | 4,900,0000 |
| Google Jam | 16,000 | 55,550 | 10,965,120 |
| Malware | 1500 | 360,000 | 45,650,214 |
| Total | 21,050 | 428,460 | 61,515,334 |

**Table 3**
F-result.

| Characteristics | Prec. | Rec. | $F_1$ |
|---|---|---|---|
| Structure | 0.95 | 0.87 | 0.91 |
| Optimization | 0.84 | 0.79 | 0.81 |
| Knowledge | 0.92 | 0.88 | 0.90 |
| Expertise | 0.97 | 0.9 | 0.93 |
| Overall | 0.89 | 0.86 | 0.87 |
| **BinChar** | **0.98** | **0.91** | **0.94** |

highest rate and we believe the reason behind this is that each programmer should follow the standard coding instructions which restrict him/her to have their own code traits.

*Authorship clusters*

For the verification of how effective *BinChar* is, we tested against a cluster of real-world malware data sets. *BinChar* has been applied to nearly 500 GB of samples of malware supplied by the Kaggle Microsoft Malware Classification Challenge (BIG, 2015) (Big Game Hunting, 2015), the aim of which is the prediction of malware family classes. Because of the way malware operates, the authors' identities are not known. It is recognized that attack operations are difficult to attribute and do not often come to court, even when there is clear evidence present (Big Game Hunting, 2015; Techniqal report and Resource, 2012). This makes it extremely difficult to establish a level of truth that can be regarded as entirely reliable. Because of this, *BinChar* regards each family as having an author(s) that share a number of characteristics. The statistical basis is outlined in Table 4. The functions are labeled by employing extant tools: IDA and pro and Paradyn are used to label library-related functions, and the BinComp (Rahimian et al., 2015) tool is used to label functions related to compilers. The samples have different varieties of the same basic malware and so we can make the assumption that each variation displaying similar characteristics was created by the same author(s). The BinComp tool (Rahimian et al., 2015) is used to obtain the number of compiler functions, and the fifth column illustrates the number of library functions obtained using F.L.I.R.T technology (HexRays, 2011) and Paradyn. Looking at Table 4, it is clear that there is quite a high percentage of compiler functions, e.g. the percentage of compiler functions in the Lollipop family is 30%.

**Cluster Quality**. One of the greatest challenges was to calculate how reliable the results obtained by the use of a clustering algorithm were (Bayer et al., 2009). To make this calculation, we can either calculate how many clusters there are, how many samples each cluster has on average, and relative sums for all pairwise distances for a cluster; alternatively a few clusters can be selected and the fact that there are similar samples in these clusters can be manually verified. Given that we had access to used samples we verified the correctness of the clusters manually.

**Table 4**
Characteristics of malware datasets (AF): Assembly Functions (CF): Compiler Functions (LF): Library Function.

| Malware | # of variants | # of BF | # of CF | # of LF |
|---|---|---|---|---|
| Ramnit | 4 | 5285 | 1601 | 50 |
| Lollipop | 3 | 3510 | 1054 | 100 |
| Kelihos | 2 | 1924 | 847 | 74 |
| Vundo | 4 | 7923 | 2410 | 219 |
| Simda | 2 | 2100 | 689 | 105 |
| Tracur | 2 | 1657 | 787 | 100 |
| Obfuscator.ACY | 3 | 2762 | 986 | 310 |
| Gatak | 2 | 2054 | 860 | 174 |

**Table 5**
Clustering results based on the features used in existing systems (TC): the total number of clusters (CC): the percentage of correct clusters (WC): the percentage of wrong clusters.

| Malware | Optimization | | | Structural | | | Knowledge | | | Expertise | | | Overall | | | BinChar | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | TC | CC | WC | TC | CC | WC | TC | CC | WC | TC | CC | WC | TC | CC | WC | TC | CC | WC |
| Ramnit | 145 | 60% | 13% | 110 | 47% | 25% | 208 | 58% | 17% | 150 | 61% | 12% | 210 | 75% | 22% | 345 | 82% | 10% |
| Lollipop | 90 | 75% | 18% | 185 | 59% | 28% | 220 | 72% | 21% | 178 | 74% | 25% | 198 | 68% | 20% | 295 | 78% | 18% |
| Kelihos | 41 | 88% | 8% | 17 | 90% | 4% | 75 | 54% | 25% | 80 | 75% | 19% | 76 | 80% | 17% | 145 | 89% | 14% |
| Vundo | 200 | 62% | 14% | 89 | 48% | 38% | 384 | 79% | 24% | 215 | 82% | 18% | 289 | 72% | 27% | 544 | 90% | 7% |
| Simda | 52 | 64% | 21% | 41 | 92% | 5% | 109 | 82% | 19% | 77 | 81% | 14% | 100 | 67% | 24% | 124 | 84% | 15% |
| Tracur | 44 | 89% | 9% | 53 | 83% | 12% | 124 | 51% | 40% | 68 | 79% | 14% | 80 | 87% | 11% | 105 | 82% | 10% |
| Obfuscator.ACY | 30 | 78% | 21% | 45 | 74% | 24% | 89 | 89% | 11% | 114 | 80% | 17% | 97 | 75% | 15% | 109 | 84% | 14% |
| Gatak | 29 | 67% | 20% | 51 | 87% | 12% | 79 | 78% | 16% | 105 | 84% | 12% | 98 | 80% | 17% | 127 | 89% | 10% |

**Table 6**
Subset of characteristics found by *BinChar* in certain malware families (✓) means the characteristic is found in Mirai botnet and it is found in that malware family, whereas (✗) means the characteristic is found in Mirai botnet and it was not found in that malware family.

| Characteristics | Ramnit | Encpk | Kelihos | Zaccess | Vobfus | Sality | Gatak | Asterope | Turla |
|---|---|---|---|---|---|---|---|---|---|
| Error messages | ✓ | ✗ | ✓ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ |
| Control commands | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ |
| Preference in keywords | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✓ | ✗ | ✗ |
| Reusing free code | ✗ | ✗ | ✓ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ |
| Using specific API | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Using specific encryption | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✓ |
| Modifying constants | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ |
| Dynamic memory allocation | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ |
| Runtime protection strategies | ✗ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ | ✓ |
| Compiler security features | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ | | ✓ | ✗ |
| Files vs memory | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Global variables vs. local variables | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Accessing closed files | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Accessing freed memory | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ |
| Function overloading | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ |
| Parameter ordering | ✗ | ✗ | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ |

*The correctness of clusters*

For the assessment of the quality of the results for the cluster algorithm, two metrics were introduced, wrong clusters (WC) and correct clusters (CC). The aim of CC is to discover how accurate the clustering algorithm is at making distinctions between samples with varied author characteristics. WC aims to measure percentages for custom functions of various classes in the same cluster.

Following on from this, we applied our tool to cluster functions in accordance with similarities in authorship characteristics, employing the standard k-mean as previously described. We then undertook manual analysis of the clusters obtained so that they could be classified as either CC or WC, as illustrated by Table 5. By analyzing the clusters we obtained, we discovered a different group of characteristics, which included means of configuration, habits of memory allocation, employment of security rules, employment of variables, etc.

## Characterizing the *Mirai* Botnet

One challenge in applying *BinChar* to real world malware is the lack of ground truth concerning the attribution of authorship due to the nature of malware. Also, whether a malware package is created by an individual or an organization is generally unconfirmed. Those limitations partially explain the fact that few research efforts have been seen on this subject. We present in this section a case study by applying *BinChar* to Mirai botent and compare the extracted characteristics with 360,000 samples collected from various databases where 40% of these samples were packed and obfuscated and we unpacked them in our security laboratory. Some of our samples include Bunny, Babar, Stuxnet, Flame, Duqu, Zeus, Citadel, zaccess, encpk, sality, etc.

**An overview**. Mirai is a DDoS botnet that has discovered by

MalwareMustDie team in August 2016. It has been considered as one of the biggest DDoS attacks on Internet which causes to shut down a major parts of Internet. It has been created using ELF binaries. The statistics of Mirai binaries is introduced in Table 7. As shown in Table 7, the samples have different platform architectures such as MIPS and PowerPc. *BinChar* can handle multiple architectures since it uses some features (Optimization characteristics and Overall characteristics) that are cross-architecture independent.

**Findings**. Our tool is able to find some authorship characteristics links between Mirai botnet and very few sample from our dataset. *BinChar* extracts various characteristics that are related to the author(s) of Mirai botnet. It finds the following: the use of all capital letters for *config* in XML; the use of a common approach to managing functions; the characteristics of opening and terminating processes; the passing of primitive types by value, but the passing of objects by reference; the use of network resources rather than file resources; creating configurations using mostly config files; and the use of semaphores and locks. More specifically, we have found different set of characteristics.

Table 6 presents examples of authorship characteristics. It can be seen that certain characteristics, such as expertise characterises

**Table 7**
Statistics about the Mirai binaries used in the case study.

| Sample | Size (KB) | No. of binary functions | Platform |
|---|---|---|---|
| File1.ARM.ELF | 1.14 | 12 | ARM |
| Mirai.arm | 5 | 79 | ARM |
| File2.MIPS.ELF | 4 | 19 | MIPS |
| mirai.sh4 | 3 | 11 | SuperH |
| mirai.sparc | 3 | 13 | SPARC |
| mirai.x86 | 7 | 31 | x86 |
| mirai.ppc | 3 | 14 | PowerPC |

**Table 8**
The number of extracted features from Mirari botnet (●) means that cannot be extracted.

| Sample | Characteristics | | | | |
|---|---|---|---|---|---|
| | Structure | Optimization | Knowledge | Expertise | Overall |
| File1.ARM.ELF | 450 | 180 | ● | ● | 3800 |
| Mirai.arm | 721 | 225 | ● | ● | 3400 |
| File2.MIPS.ELF | 400 | 190 | ● | ● | 2900 |
| mirai.x86 | 950 | 428 | 215 | 147 | 4500 |
| mirai.ppc | 145 | 175 | ● | ● | 3010 |
| mirai.sparc | 200 | 215 | ● | ● | 3700 |
| mirai.sh4 | 315 | 160 | ● | ● | 3150 |

**Table 9**
Similarity between authorship characteristics found in Mirai botnent and characteristics found in malware binaries.

| | Bunny | Babar | Stuxnet | Flame | Zeus | Citadel | Mirai |
|---|---|---|---|---|---|---|---|
| Bunny | — | 61% | 14% | 8% | 11% | 13% | 1.2% |
| Babar | 61% | — | 9% | 10% | 2% | 8% | 2.09% |
| Stuxnet | 14% | 9% | — | 74% | 19% | 10% | 0.96% |
| Flame | 8% | 10% | 74% | — | 14% | 6% | 0.71% |
| Zeus | 11% | 2% | 19% | 14% | — | 80% | 0 |
| Citadel | 13% | 8% | 10% | 6% | 80% | — | 0.09% |
| Mirai | 1.2% | 2.09% | 0.96% | 0.71% | 0 | 0.09% | — |

including control command strings and dynamic memory allocation. Also, some examples on knowledges characteristics including the use of API calls, runtime protection strategies, etc. We have found another characteristics that belong to code traits such as accessing closed files or modifying constants. We studied the number of features extracted from Mirai botnet samples (Table 8). For example, the number of features that are related to structure characteristics authors is 175 extracted from mirai.ppc. As shown, there are some features could not be extracted from files that are not x86-based architecture. We leave this issue for future work venue. We have studied the impact of each authorship characteristics on accuracy (Fig. 3). For example, the $F_1$ score is 0.7 between Mirai.x86 and Ramnit when we compare them based on the knowledge characteristics. While the average $F_1$ score is 0.42 between Ramnit and Mirai.x86. As shown, the highest $F_1$ score is 0.83 between Mirai.x86 and Vobfus when we compare them based on optimization characteristics.

**Measuring the degree of similarity between ground truth datasets and malware binaries.** As another verification of the correctness of the findings, we measured the degree of similarity between the Mirai Botnet here and other datasets for which we have the ground truth (e.g., Google code jam) to see how likely such a degree of similarity could come from shared authorship characteristics. The goal of computing the degree of similarity is to determine whether the authorship characteristics found in the Mirai botnet are present to the same degree in conventional binaries, which will reveal whether these characteristics are indeed specific to malware writers. To provide an even more convincing verification, we computed the similarity scores between related Mirai botnet samples and the rest of the available dataset. BinChar found a similarity of 7% with those characteristics in the Google code jam dataset and 17% with those characteristics in the Github dataset. We believe that one of the main reasons for the high similarity is that the programmers participating in the Github may have greater expertise, more extensive background knowledge, and better skills than the programmers who participate in Google code jam.
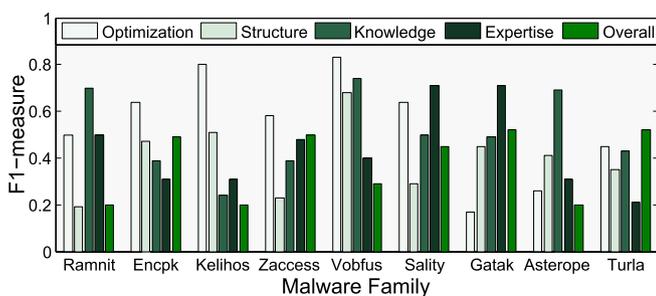
**Measuring similarity between authorship characteristics in malware binaries.** In this section, the goal is to assess the similarity between malware binaries by reporting the similarity in terms of authorship characteristics (Table 9).

We compare similarity between the author characteristics extracted from Mirai botnet to different sets of real malware: i) Bunny and Babar; ii) Stuxnet and Flame; and iii) Zeus and Citadel. These malware are selected based on researchers' claims that each of these pairs of malware originates from the same set of authors (Techniqal report and Mcafee, 2011; Techniqal report and Resource, 2012; Big Game Hunting, 2015). We observed that the similarity between the authorship characteristics found in Mirai botnent and characteristics found in malware binaries is very low. For instance, the similarity between Mirai and Zeus is 0. In the meantime, our tool is able to find common authorship characteristics among other malware families. For example, the similarity between Bunny and Babar is about 60%. This finding supports the claim that this pair originates from the same set of authors (Big Game Hunting, 2015).

### Limitations and concluding remarks

Our work has certain limitations. First, the system is unlikely to remain accurate if the authors used advanced obfuscation techniques to evade detection. Second, although we have tested our work on relatively large datasets, our dataset can still be enriched in terms of both scale and scope. Third, the features used by BinChar are static-based; as such, BinChar cannot detect characteristics that require dynamic features. Fourth, there is a room for performance improvement by including some techniques such as Map-reduced methods. Finally, since BinChar currently supports the x86 ISA, other ISAs such as ARM should also be considered. Our future work aims to extend BinChar to tackle these limitations.

To conclude, we have presented the first known effort on characterizing the author of binary code based on personnel characteristics. Previous existing works have only employed artificial datasets, whereas we included more realistic datasets. We also applied our system to known malware. In summary, our system demonstrates superior results on more realistic datasets and real malware and can detect the presence of multiple authors.

### Acknowledgments

### References

Alrabaee, S., Saleem, N., Preda, S., Wang, L., Debbabi, M., 2014. Oba2: an onion approach to binary code authorship attribution. Digit. Invest. 11, S94–S103.
Alrabaee, S., Shirani, P., Wang, L., Debbabi, M., Hanna, A., 2018. On leveraging coding



**Fig. 3.** The F1 score between Mirai.x86 and the malware families in Table 6.

habits for effective binary authorship attribution. In: European Symposium on Research in Computer Security. Springer, pp. 26–47.

Bayer, U., Comparetti, P.M., Hlauschek, C., Kruegel, C., Kirda, E., 2009. Scalable, behavior-based malware clustering. NDSS 9, 8–11.

Big Game Hunting: Nation-State Malware Research, 2015. BlackHat. https://www.blackhat.com/docs/us-15/materials/us-15-MarquisBoire-Big-Game-Hunting-The-Peculiarities-Of-Nation-State-Malware-Research.pdf.

A. Caliskan-Islam, F. Yamaguchi, E. Dauber, R. Harang, K. Rieck, R. Greenstadt, A. Narayanan, When Coding Style Survives Compilation: De-anonymizing Programmers from Executable Binaries, arXiv preprint arXiv:1512.08546.

Caliskan-Islam, A., Harang, R., Liu, A., Narayanan, A., Voss, C., Yamaguchi, F., Greenstadt, R., 2015. De-anonymizing Programmers via Code Stylometry. USENIX.

Citizen Lab, 2015. university of Toronto, Canada. https://citizenlab.org/.

Farnstrom, F., Lewis, J., Elkan, C., 2000. Scalability for clustering algorithms revisited. ACM SIGKDD Explorations Newsletter 2 (1), 51–57.

PEfile, 2012 accessed on Nov, 2016. http://code.google.com/p/pefile/.

Frantzeskou, G., 2004. Source Code Authorship Analysis for Supporting the Cybercrime Investigation Process, pp. 470–495.

Glorot, X., Bordes, A., Bengio, Y., 2011. Deep sparse rectifier neural networks. In: Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics, pp. 315–323.

HexRays: IDA Pro, 2011. https://www.hex-rays.com/products/ida/index.shtml.

G. Huang, Z. Liu, K. Q. Weinberger, L. van der Maaten, Densely Connected Convolutional Networks, arXiv preprint arXiv:1608.06993.

Huang, H., Youssef, A.M., Debbabi, M., 2017. Binsequence: fast, accurate and scalable binary code reuse detection. In: Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security. ACM, pp. 155–166.

Karbab, E.B., Debbabi, M., Alrabaee, S., Mouheb, D., 2016. Dysign: dynamic fingerprinting for the automatic detection of android malware. In: Malicious and Unwanted Software (MALWARE), 2016 11th International Conference on. IEEE, pp. 1–8.

Kirat, D., Nataraj, L., Vigna, G., Manjunath, B., 2013. Sigmal: a static signal processing based malware triage. In: Proceedings of the 29th Annual Computer Security Applications Conference. ACM, pp. 89–98.

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

lab, L., 2014. Deep Learning Tutorial Release 0.1. university of Montreal. http://deeplearning.net/tutorial/deeplearning.pdf.

Maas, A.L., Hannun, A.Y., Ng, A.Y., 2013. Rectifier Nonlinearities Improve Neural Network Acoustic Models. ICML.

McCabe, T.J., 1976. A complexity measure. IEEE Trans. Softw. Eng. SE-2 (4), 308–320.

Meng, X., 2016. Fine-grained binary code authorship identification. In: Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering. ACM, pp. 1097–1099.

Moran, N., Bennett, J., 2013. Supply Chain Analysis: from Quartermaster to SunShop, vol. 11. FireEye Labs.

Nagano, Y., Uda, R., 2017. Static analysis with paragraph vector for malware detection. In: Proceedings of the 11th International Conference on Ubiquitous Information Management and Communication. ACM, p. 80.

Nataraj, L., Karthikeyan, S., Jacob, G., Manjunath, B., 2011. Malware images: visualization and automatic classification. In: Proceedings of the 8th International Symposium on Visualization for Cyber Security. ACM, p. 4.

Nataraj, L., Kirat, D., Manjunath, B., Vigna, G., 2013. Sarvam: search and retrieval of malware. In: Proceedings of the Annual Computer Security Conference (ACSAC) Worshop on Next Generation Malware Attacks and Defense. NGMAD.

Rahimian, A., Shirani, P., Alrbaee, S., Wang, L., Debbabi, M., 2015. Bincomp: a stratified approach to compiler provenance attribution. Digit. Invest. 14, S146–S155.

Rosenblum, N., Zhu, X., Miller, B.P., 2011. Who wrote this code? identifying the authors of program binaries. In: Computer Security–ESORICS 2011. Springer, pp. 172–189.

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

Schkufza, E., 2015. Stochastic Program Optimization for X86 64 Binaries. Ph.D. thesis. STANFORD UNIVERSITY.

Sharma, R., 2016. Data-driven Verification. Ph.D. thesis. Stanford University.

Shin, E.C.R., Song, D., Moazzezi, R., 2015. Recognizing Functions in Binaries with Neural Networks. USENIX.

Tate, R., Stepp, M., Tatlock, Z., Lerner, S., 2009. Equality saturation: a new approach to optimization. In: ACM SIGPLAN Notices, vol. 44. ACM, pp. 264–276.

Taylor, Q.C., Stevenson, J.E., Delorey, D.P., Knutson, C.D., 2008. Author entropy: a metric for characterization of software authorship patterns. In: Third International Workshop on Public Data about Software Development ((WoPDaSD08), p. 6.

Techniqal Report, Mcafee, 2011. www.mcafee.com/ca/resources/wp-citadel-trojan-summary.pdf.

reportTechniqal report, Resource 207: Kaspersky Lab Research proves that Stuxnet and Flame developers are connected, http://www.kaspersky.com/about/news/virus/2012/.

The GitHub Repository, 2016. https://github.com/.

The Google Code Jam, 2008-2015. http://code.google.com/codejam/.

Tristan, J.-B., Govereau, P., Morrisett, G., 2011. Evaluating value-graph translation validation for llvm. ACM Sigplan Not. 46 (6), 295–305.

Wagner, M., Rind, A., Thür, N., Aigner, W., 2017. A knowledge-assisted visual malware analysis system: design, validation, and reflection of kamas. Comput. Secur. 67, 1–15.

Wei, Y., Xia, W., Lin, M., Huang, J., Ni, B., Dong, J., Zhao, Y., Yan, S., 2016. Hcp: a flexible cnn framework for multi-label image classification. IEEE Trans. Pattern Anal. Mach. Intell. 38 (9), 1901–1907.

Woldring, D.R., Holec, P.V., Hackel, B.J., 2016. Scaffoldseq: software for characterization of directed evolution populations. Proteins: Structure, Function, and Bioinformatics 84 (7), 869–874.