

Random CapsNet Forest Model for Imbalanced Malware Type Classification Task

Aykut Çayır*, Uğur Ünal, Hasan Dağ

Management Information Systems Department, T. C. Kadir Has University, Istanbul, Turkey

Abstract

Behavior of malware varies concerning the malware types, which affects the strategies of the system protection software. Many malware classification models, empowered by machine and/or deep learning, achieve superior accuracies for predicting malware types. Machine learning-based models need to do heavy feature engineering work, which affects the performance of the models greatly. On the other hand, deep learning-based models require less effort in feature engineering when compared to that of the machine learning-based models. However, traditional deep learning architectures components, such as max and average pooling, cause architecture to be more complex and the models to be more sensitive to data. The capsule network architectures, on the other hand, reduce the aforementioned complexities by eliminating the pooling components. Additionally, capsule network architectures based models are less sensitive to data, unlike the classical convolutional neural network architectures. This paper proposes an ensemble capsule network model based on the bootstrap aggregating technique. The proposed method is tested on two widely used, highly imbalanced datasets (Maling and BIG2015), for which the-state-of-the-art results are well-known and can be used for comparison purposes. RCNF achieves the highest F-Score, which is 0.9820, for the BIG2015 dataset and F-Score, which is 0.9661, for the Maling dataset. RCNF reaches the-state-of-the-art with fewer trainable parameters than other competitors.

*Corresponding author

Email address: aykut.cayir@khas.edu.tr (Aykut Çayır)

Keywords: Capsule networks, Malware, Ensemble model, Deep learning, Machine learning

1. Introduction

Malware type classification is as important as malware detection problem because system protection software makes their strategies concerning malware family types. Malware families have different behaviors and effects on a computer system. Each malware family uses different resources, files, ports, and other components of operating systems. For example, malware in the online banking systems aim to perform fraud, steal private information of users, and use different spreading behaviors [1, 2]. In addition to this, due to the trends in technology, new malware types occur almost daily. Thus, most of the computers, smartphones, and other digital systems are vulnerable to new malware. In this case, many zero-day attacks are performed [3]. The raising of the number of malware makes using the big data techniques crucial for malware analysis [4].

Malware type classification is the most common problem in the cybersecurity domain, because strategies of protection systems vary with respect to malware family type. Malware type classification problem is broadly dealt with in three different ways: static, dynamic, and image-based [5, 6, 7]. This paper focuses on the image-based malware family classification problem. However, malware family type classification is an imbalanced task, so this makes many models unsuccessful at predicting rare classes. To this end, two imbalanced datasets are used and the results are compared to other models in the literature.

This paper proposes a new model named **Random Capsule Network Forest (RCNF)** based on bootstrap aggregation (bagging) ensemble technique and capsule network (CapsNet) [8, 9]. The main motive behind the proposed method is to reduce the variance of different CapsNet models (as weak learners) using bagging. In this perspective, the main contributions of this paper can be listed as follows:

- The paper introduces the first application of CapsNet in the field of malware type classification. Although image-based malware classification is a broad research and application area, there is no research and application of CapsNet in the literature in our best knowledge.

- The paper uses the first ensemble model of CapsNet. The key idea of creating an ensemble of CapsNet is assuming a single CapsNet model as a weak classifier like a decision tree model. In this way, an ensemble model of CapsNet can be easily created using bootstrap aggregating. The assumption that CapsNet is a weak learner increases the performance of a single CapsNet for two different well-known malware datasets, which are highly imbalanced.
- The proposed model uses simple architecture engineering instead of complex convolutional neural network architectures and domain-specific feature engineering techniques. In addition to this, CapsNet does not require to use transfer learning, and the model is easily trained from scratch. Because of that, the created network and its ensemble version have reasonably a lower the number of parameters.
- The proposed model is compared with the latest studies that use deep neural networks for image-based malware classification tasks. For a fair comparison, especially, the last studies using the Maling and the BIG2015 datasets are chosen and compared with the proposed method.

Image-based malware classification is a broad research and application area. At the same time, deep learning drives the computer vision and image processing researches. Many deep convolutional neural networks have proven their success in image processing. CapsNet is the most important deep convolutional neural architecture that removes pooling to avoid losing the spatial features of images. This is the power of CapsNet comparing to classical CNNs. Therefore, the number of applications of CapsNet is increasing in image processing. The main motivation of this paper is to design a simple and accurate classifier for imbalanced malware type classification problem using bagging [8, 10] and CapsNet architecture. This paper also presents a detailed comparison of the proposed model with the other models in the literature.

The paper is organized as follows. Section 2 presents a literature survey for CapsNet applications and previous malware analysis studies. The methodology of the paper is described in Section 3, whereas Section 4 gives details of the inspiring model and the proposed model. In Section 5, the test results are discussed and many comparisons with related works

published in the last years are listed; and finally, Section 6 provides the concluding remarks.

2. Related Work

There are many different ways to represent malware files for machine learning-based identification. One of them is to extract features from the application programming interface (API) calls of malware. For example, Alazab [11] proposed a framework to get features statically and dynamically from malware API calls. He used similarity mining and machine learning to profile and classify malware. He obtained a 0.966 received-operating-curve (ROC) score in the malware dataset containing 66,703 samples (Malign or Benign) with the k-nearest neighbor algorithm. Moreover, Azab et al. [12] focused on grouping malware in the same variants using hashing techniques for malware binaries. They used two different Zeus datasets. The first dataset contained 856 binaries, and the second dataset contained 22 binaries. Each binary had SHA256 value. They achieved 0.999 F-Score using the k-nearest neighbors and SDHASH.

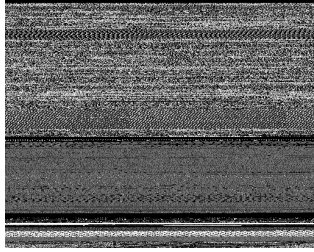
The second efficient way to feed machine learning algorithms to classify malware files is image-based representation. In this work, we have focused on image-based malware type classification. For example, Nataraj et al. [13] converted malware files to greyscale images to represent malware. They extracted GIST features from malware images, and then they classified malware family types, using the euclidean k-nearest neighbors algorithm. They reached 0.98 classification accuracy for the dataset that has 9339 samples and 25 malware families. Similarly, Kancherla et al. [14] used image-based malware representation to feed the support vector machine to classify malware files malignant or benign. They extracted three different features, such as intensity-based, wavelet-based, and Gabor based. Their dataset contained 15,000 malignant and 12,000 benign samples. They attained 0.979 ROC score. These studies utilized traditional machine learning algorithms, such as k-nearest neighbors, and the support vector machines. These algorithms required to extract good features from images to classify malware types with a high performance. After an impactful success of a deep convolutional neural network (CNN) on the Imagenet dataset, a new era started in computer vision [15]. CNNs could classify images using raw pixel values without complex feature engineering methods.

Image-based malware classification is a broad research area, which is affected by deep convolutional networks. In this regard, one of the most important applications of CNNs is the transfer learning, which is useful and successful for the balanced and relatively small size datasets [16]. Sang Ni et al. [7] created greyscale image files using SimHash bits of malware. They obtained 0.9926 accuracy on the dataset containing 10,805 samples, using a CNN classifier. There are many deep learning models to classify malware types, but we used some of them in the experiment part of the paper.

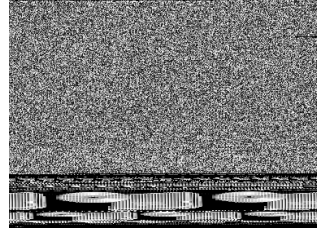
CapsNet, as a new CNN structure, has been implemented in 2017 [9], especially in the health domain [17], CapsNet have many applications in the literature. For instance, Afshar et al. [18] use CapsNet for brain tumor classification problems like classification of breast cancer histology images of [19]. Mobiny et al. [20] create a fast CapsNet architecture for lung cancer diagnosis. Another important application area of CapsNet is object segmentation. LaLonde et al. [21] use CapsNet for object segmentation. Traditional CNN structures, on the other hand, are used in **Generative Adversarial Networks**, GANs. CapsNet is very useful to make GANs better by removing the weakest point of these CNNs [22]. The snippet studies above show that CapsNet is a promising architecture against the standard CNN. In the literature, although there are many applications of CapsNets, there is a missing and important area. This area is a computer and information security. This gap can be seen easily in the pre-print version of a survey about CapsNets [23].

Another crucial issue in malware classification is imbalanced datasets. Eбенуwa et al. [24] pointed imbalanced classification problem in binary classification. They inspected three different techniques, such as sampling-based, algorithm modifications, and cost-sensitive approaches. They proposed variance ranking feature selection techniques to get better results in imbalanced datasets for binary classification problems.

To this end, this paper aims to develop a malware classification model based on an ensemble of CapsNet architecture for imbalanced malware datasets, which is the first application of CapsNets in the malware classification domain.



(a) An image example from Family Adier.



(b) An image example from Family Fakerean.

Figure 1: Malware image samples obtained from byte files using the algorithm described in [13].

3. Methodology

3.1. Malware Datasets

There are many open research issues in malware classification. These issues can be listed such as class imbalance, concept drift, adversarial learning, interpretability (explainability) of the models, and public benchmarks [25]. In this paper, our model called RCNF focuses on the class imbalance issue. Thus, the base CapsNet and the proposed RCNF models have been tested on two very well-known malware datasets called Malimg and Microsoft Malware 2015 (BIG2015). These datasets are highly imbalanced in terms of the number of classes. This section describes these datasets.

3.1.1. Malimg

Nataraj et al. introduced a new malware family type classification approach based on visual analysis, converted binaries into greyscale images and they published these images as a new malware dataset called Malimg [13]. This dataset has 9339 samples and 25 different classes. Table 1 presents the number of samples for each malware family. This distribution shows that the dataset is highly imbalanced.

Fig. 1 shows the malware images created from the byte files. All images are single-channel and are resized to 224×224 for CapsNet architecture. This size is the largest value that can be processed in our computer system.

Table 1: Sample Distribution for each Malware Family.

| No. | Family Name | Number of Samples |
|-----|----------------|-------------------|
| 1 | Allaple.L | 1591 |
| 2 | Allaple.A | 2949 |
| 3 | Yuner.A | 800 |
| 4 | Lolyda.AA 1 | 213 |
| 5 | Lolyda.AA 2 | 184 |
| 6 | Lolyda.AA 3 | 123 |
| 7 | C2Lop.P | 146 |
| 8 | C2Lop.gen!g | 200 |
| 9 | Instantaccess | 431 |
| 10 | Swizzot.gen!I | 132 |
| 11 | Swizzor.gen!E | 128 |
| 12 | VB.AT | 408 |
| 13 | Fakerean | 381 |
| 14 | Alueron.gen!J | 198 |
| 15 | Malex.gen!J | 136 |
| 16 | Lolyda.AT | 159 |
| 17 | Adialer.C | 125 |
| 18 | Wintrim.BX | 97 |
| 19 | Dialplatform.B | 177 |
| 20 | Dontovo.A | 162 |
| 21 | Obfuscator.AD | 142 |
| 22 | Agent.FYI | 116 |
| 23 | Autorun.K | 106 |
| 24 | Rbot!gen | 158 |
| 25 | Skintrim.N | 80 |

3.1.2. Microsoft Malware 2015 (BIG2015)

BIG2015 dataset has been released as a Kaggle competition [26, 27]. Table 2 presents the sample distribution for each malware family in BIG-2015 dataset. The distribution shows that the dataset is highly imbalanced; and *Simda* is the toughest malware family to be predicted for the dataset. The dataset contains 10868 BYTE (bytes) files and 10868 ASM (assembly code) files and 9 different malware family types. BIG2015, unlike the Maling dataset, contains raw files. Thus, a file from the BIG2015

dataset is opened in the byte mode and then the file is read by 256 sized chunks till the end of the file. Finally, the buffer is converted to array and the array is saved as a greyscale image into the file system. All processes are described in Fig. 2. This method is the most common way to convert from malware files to images [13, 28].

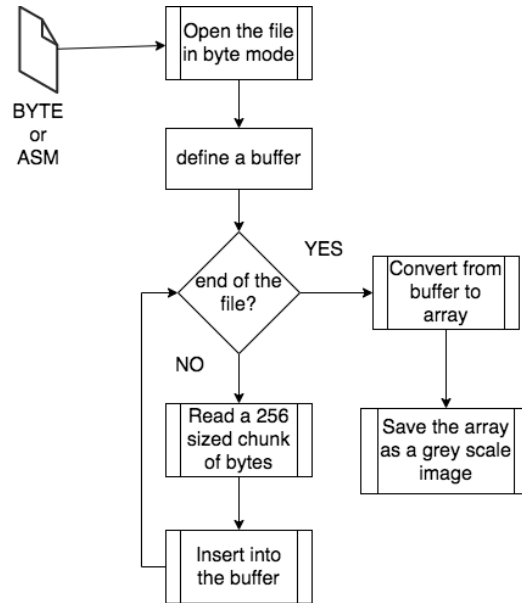


Figure 2: Flowchart of the process of a flat-file to greyscale image for the BIG2015 dataset.

Fig. 3 depicts image representations created from the BYTE and ASM files of the same malware sample in *Ramnit* malware family. All images are single channel. All images are resized 112×112 for our CapsNet architecture, because the architecture uses both BYTE and ASM image representations at the same time.

4. Model

In this section, general capsule networks, base CapsNet architecture for Malimg and base CapsNet architecture for BIG2015 are described. CapsNet architectures are different for both Malimg and BIG2015 Dataset.

Table 2: Number of Samples for Each Malware Family in BIG2015 Dataset.

| No. | Family Name | Number of Image |
|-----|----------------|-----------------|
| 1 | Ramnit | 1541 |
| 2 | Lollipop | 2478 |
| 3 | Kelihos_ver3 | 2942 |
| 4 | Vundo | 475 |
| 5 | Simda | 42 |
| 6 | Tracur | 751 |
| 7 | Kelihos_ver1 | 398 |
| 8 | Obfuscator.ACY | 1228 |
| 9 | Gatak | 1013 |

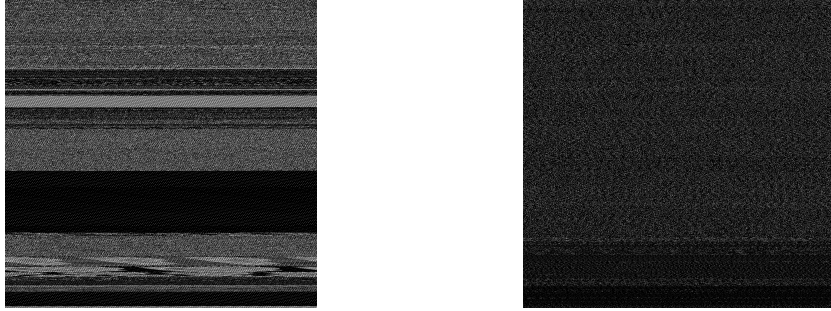
4.1. Capsule Networks

Capsule networks are special architectures of convolutional neural networks aiming to minimize information loss because of max pooling [9]. This method is the weakest point for preserving spatial information [19]. A CapsNet contains capsules similar to autoencoders [29, 9]. Each capsule learns how to represent an instance for a given class. Therefore, each capsule creates a fixed-length feature vector to be input for a classifier layer without using max pooling layers in its internal structure. In this way, this capsule structure aims to preserve texture and spatial information with minimum loss.

Sabor et al. propose an efficient method to train CapsNet architectures [9]. This method is called a dynamic routing algorithm, which uses a new non-linear activation function called squashing shown in (1). This equation emphasizes that short vectors are shrunk to almost zero and long vectors are shrunk to 1 [9]. In this equation, v_i is the output of i^{th} capsule and s_i shows the total input of this capsule.

$$v_i = \frac{\|s_i\|^2}{1 + \|s_i\|^2} \frac{s_i}{\|s_i\|} \quad (1)$$

Visualizing the squash activation function described in (1) is hard because its input is a high dimensional vector. If the activation function can be thought of as a single variable function, as described in [30], then the behavior of the function and its derivative can be visualized, as in Fig. 4.



(a) An image example obtained from
BYTE File
in Ramnit Family

(b) An image example obtained from
ASM File
in Ramnit Family

Figure 3: BIG2015 image samples from BYTE and ASM files using the algorithm described in [13].

A basic CapsNet architecture contains two parts: the standard convolution blocks and the capsule layer as shown in Fig. 5. A convolution block is made from a combination of convolution filters and ReLU activation function. At the end of the convolution block, obtained feature maps are reshaped and projected to $d - dimensional$ vector representation. This representation feeds each capsule in the capsule layer. Each capsule learns how to represent and reconstruct a given sample like an autoencoder [29] architecture. In order to learn how to reconstruct a malware sample, the capsule network will minimize reconstruction error in (2), where $x_c \in R^{d \times d}$ is the real sample in the capsule c and $\hat{x}_c \in R^{d \times d}$ is the reconstructed sample by the same capsule c . These representations are used to calculate the class probabilities for the classification task.

$$\ell_r = (x_c - \hat{x}_c)^2 \quad (2)$$

Margin loss function is used for CapsNet. This function is similar to hinge loss [31]. (3) defines the margin loss L_c for capsule c ,

$$\ell_m = y_c \times (\max(0, m - \hat{y}_c))^2 + \lambda \times (1 - y_c) \times (\max(0, \hat{y}_c - (1 - m)))^2 \quad (3)$$

where $m = 0.9$, $\lambda = 0.5$, y_c denotes actual class and \hat{y}_c represents the current prediction.

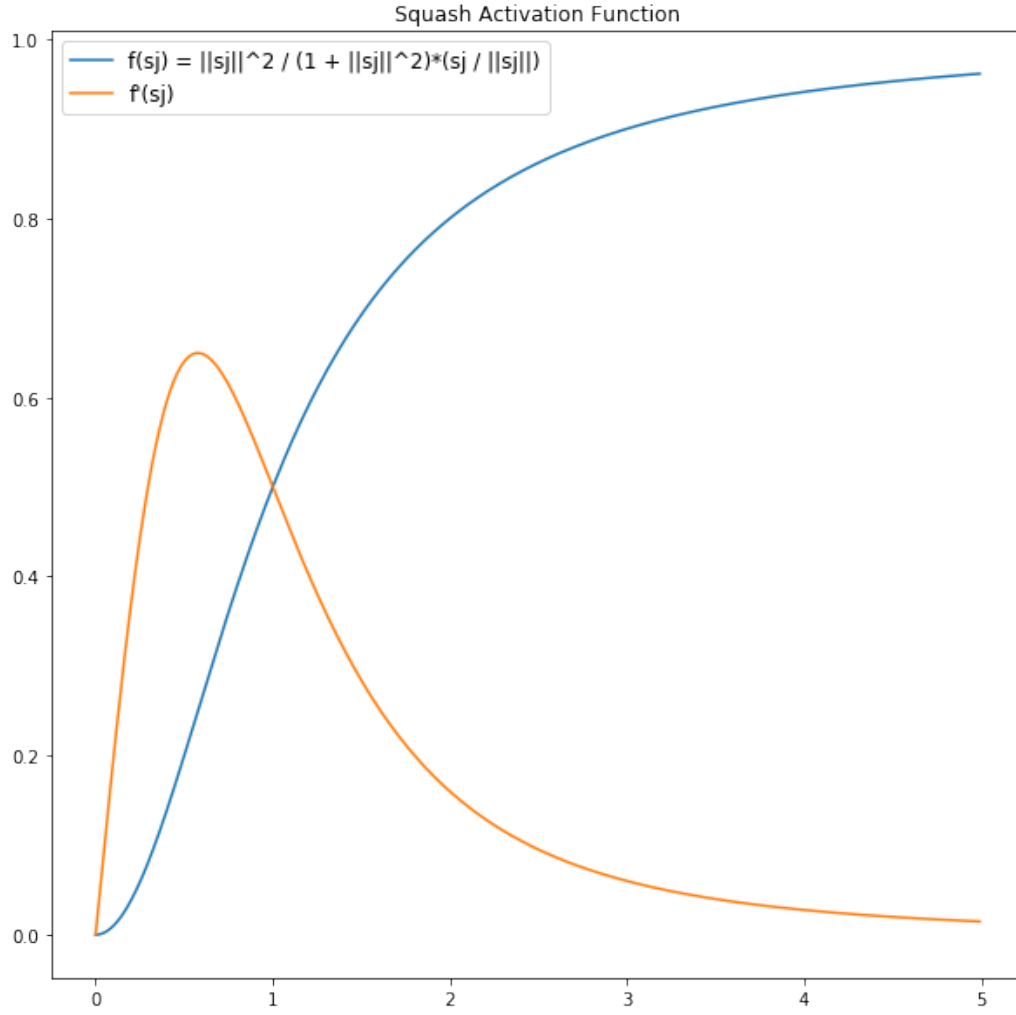


Figure 4: Squashing activation function and its derivation in 2-D plane.

$$L_c = \ell_m + 0.0005 \times \ell_r \quad (4)$$

$$L = \frac{1}{N} \sum_{n=1}^N L_c \quad (5)$$

The mean of L_c for each capsule gives the total loss in (5), where L_c is

sum of margin loss ℓ_m (as described in (3)) and reconstruction loss ℓ_r (as described in (2)). However, reconstruction loss is multiplied by 0.0005 to avoid suppressing the margin loss [9]. In order to minimize loss L , one can use the most applicable optimizer algorithm for CapsNet by Adam [32, 9]. We have observed that CapsNet cannot converge to minimum loss value with optimizers other than that of Adam. This is obviously an open issue for CapsNet studies in the future.

In image-based malware family type classification problem, there are no complex patterns that are easily detected by classical convolutional neural networks. For this reason, the predictive model must recognize the pattern of pixel distribution of the image-based malware sample. On the other hand, CapsNet can learn pixel density distribution of each malware family. Thus, a CapsNet model can be easily trained from scratch for this problem, unlike CNNs. This is the most important advantage of using CapsNet as a base classifier in our proposed model.

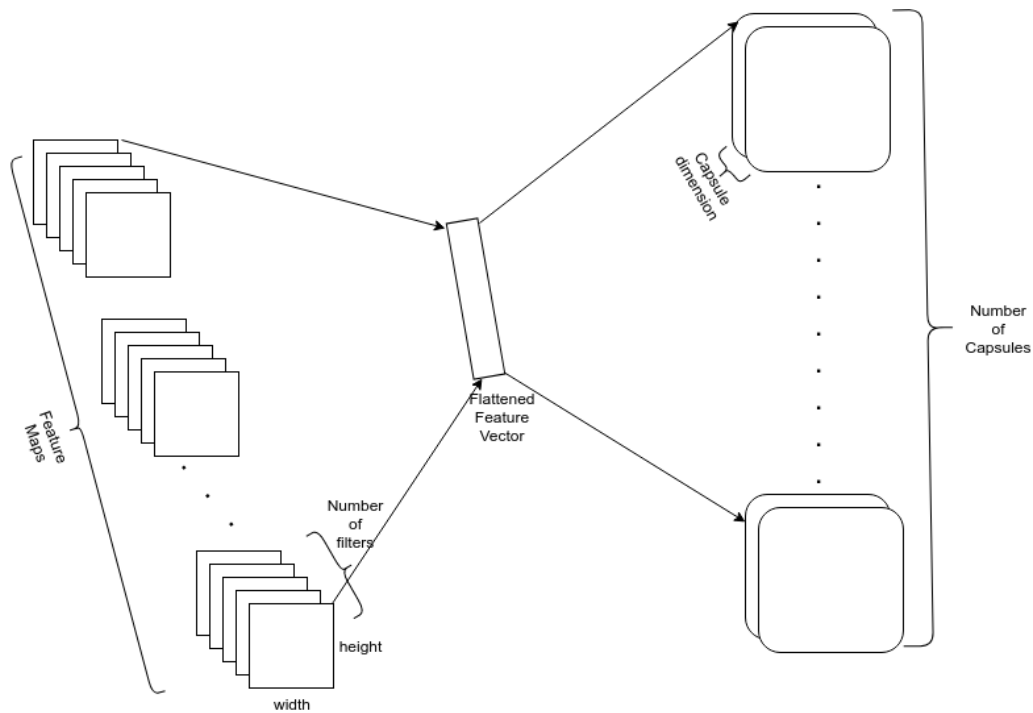


Figure 5: Basic CapsNet architecture.

Our main assumption is that CapsNet architecture will be able to successfully classify malware family types using raw pixel values obtained from malware binary and assembly files. In addition to the main assumption, this paper aims to increase CapsNet malware type classification architecture accuracy with the bagging ensemble method.

4.2. Base Capsule Network Model for Malimg Dataset

Before creating an ensemble CapsNet model, the base CapsNet estimator must be built. This architecture depends on the dataset. Thus, base CapsNet estimator architecture has a single convolution line, as shown in Fig. 6. The convolutional line contains two sequential blocks; and each block contains two sequential convolutions and ReLU layers. The first two convolutional layers have 3×3 kernels and 32 filters. The second two convolutional layers have 3×3 kernels and 64 filters. Feature maps are reshaped to 128 feature vectors. At the end of the reshape step, there is a capsule layer containing 25 capsules; the dimension of each capsule is 8; and the routing iteration is 3 of the capsule layer. This is the optimal CapsNet architecture for the Malimg dataset depending on our experiments.

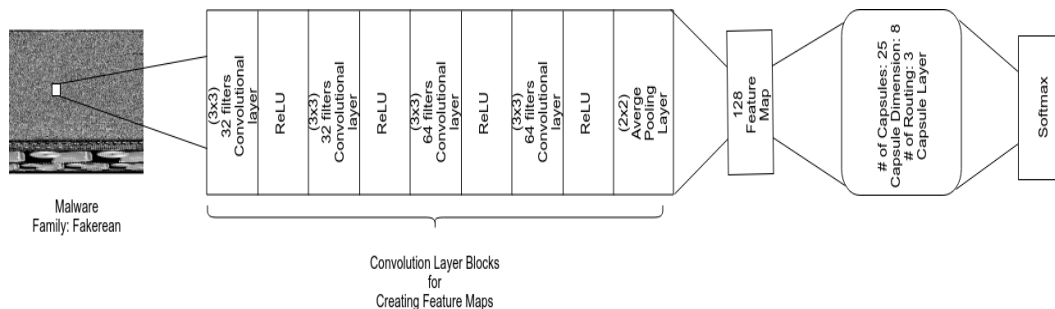


Figure 6: CapsNet architecture for Malimg dataset.

4.3. Base Capsule Network Model for BIG2015 Dataset

The BIG2015 dataset has two different files for each sample. One of them is a binary file; and the other is an assembly file. Thus, it is possible to design a CapsNet, which can be fed by two different image inputs at

the same time. Fig. 7 shows a CapsNet architecture, which has two exactly identical convolution lines. In this architecture, the first two sequential layers contain 3×3 kernels and 64 filters. The second two sequential layers contain 3×3 kernels and 128 filters.

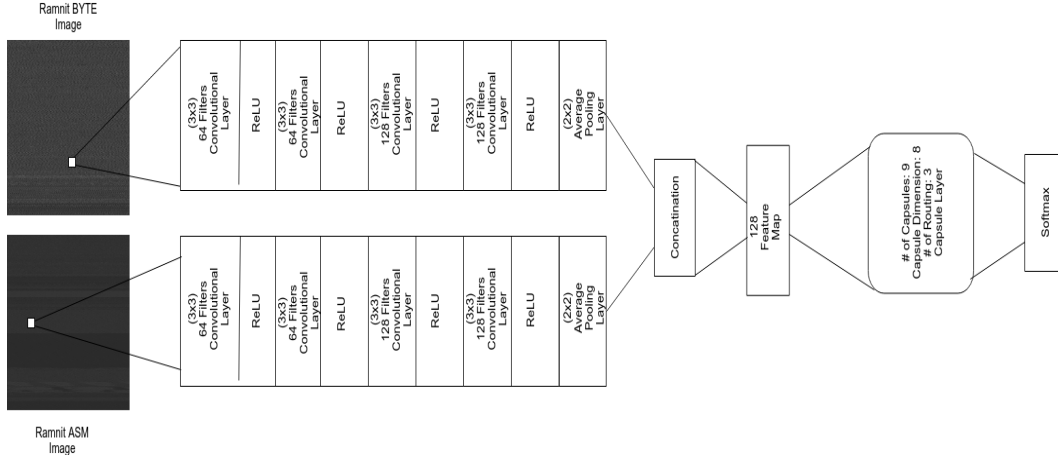


Figure 7: CapsNet architecture for BIG2015 dataset.

Features extracted from the ASM and BYTE images are concatenated; and the final feature vector is reshaped to a vector with length 128. For the next level, as an input, this feature vector feeds to a capsule layer containing 9 capsules. In this layer, the dimension of each capsule is 8; and the routing iteration is 3. This hyper-parameter set is optimal for the base CapsNet estimator for the BIG2015 dataset.

4.4. The Proposed Random CapsNet Forest Model for Imbalanced Datasets

Random CapsNet Forest (RCNF) is an ensemble model, which is inspired by the Random Forest algorithm [10]. The basic idea behind RCNF is to assume identical CapsNet models as weak learners create different training sets for each model from the original training set using the bootstrap resampling technique, as shown in Algorithm 1. The training algorithm is a variant of bootstrap aggregating (also known as bagging) [8] for CapsNet model; and bagging reduces the variance of the model while increasing robustness of the model [33]. In this paper, bagging is preferred to create an ensemble of CapsNet instead of boosting [34], because it is shown that boosting tends to overfit [35]. During the training

phase, each epoch updates the weights of the CapsNet. Therefore, the weight of the best model at the end of each epoch is saved according to the validation score to increase model performance and consistency against random weight initialization of the CapsNet.

Algorithm 1 Random CapsNet Forest Training Algorithm

```

1: procedure TRAIN(base_model, n_estimators, trainset, valset, epochs)
2:   for i ← 1, n_estimators do
3:     bs_trainset ← resample(trainset, replacement = True)
4:     for e ← 1, epochs do
5:       base_model.fit(bs_trainset)
6:       val_score ← get_accuracy(base_model, valset)
7:       if is_best_score(val_score) == True then
8:         save_weights(base_model)

```

Algorithm 2 Random CapsNet Forest Prediction Algorithm

```

1: procedure PREDICT(n_estimators, testset, numclasses)      ▷ Average
   Ensembling
2:   total_preds ← zeros_like(testset.shape[0], numclasses)
3:   for i ← 1, n_estimators do
4:     modeli ← load_model_weights(i)
5:     total_preds+ = modeli.predict(testset)
6:   preds ← total_preds/n_estimators
7:   return argmax(preds)  ▷ The final predictions of CapsNet models

```

The prediction method is described in Algorithm 2. Each weight of CapsNet model is loaded; and test samples are predicted by the model. Cumulative predicted probabilities are added onto *total_preds* variable; and this step is known as average ensembling step. At the end of the estimation loop, the index of the highest probabilities is assigned as a predicted class.

There are several limitations to the RCNF model. The first limitation is the number estimators in the RCNF model. In this implementation, an RCNF model can contain up to 10 CapsNets because of an increasing number of trainable parameters. The second limitation is the training time. Training of an RCNF with 10 CapsNets for the BIG2015 dataset

takes five hours (100 epochs for each CapsNet). This training time of the RCNF with 5 CapsNets model for the Maling dataset is decreasing (100 epochs for each CapsNet). On the other hand, the RCNF can be easily parallelized to increase efficiency in the training phase. Each CapsNet can be trained on multiple GPUs. We will develop a distributed multi-GPU version of the RCNF as a future work.

We implemented the RCNF using Tensorflow (version 1.5) [36] and Keras [37], Sklearn [38], Numpy [39] and Pandas [40]. All scripts were written in Python3. The configuration of the computer used in this study was 12GB GPU (GeForce GTX 1080 Ti) and Intel Core i9-9900K processor with 64 GB main memory for testing.

5. Experiment and Results

CapsNet and RCNF ensemble models are tested on two different datasets called Maling and BIG2015. The Maling dataset has been divided into three parts: training, validation, and test sets. The training set has 7004 samples, the validation set has 1167 samples, and the test set has 1167 samples. BIG2015 has also been divided into three parts like the Maling dataset. In the experiments of CapsNet and RCNF ensemble model for the BIG2015 dataset, the training set has 8151 samples, the validation set has 1359 samples and the test set has 1358 samples. The first experiment is made to obtain performance results of single base CapsNet estimators for each dataset. The second experiment is about the performance of the RCNF model. Model evaluation has been done in terms of accuracy, F-Score, and the number of parameters of deep neural nets. These performance metrics are defined as follows:

$$accuracy = \frac{(TP + TN)}{(TP + TN + FP + FN)} \quad (6)$$

$$F\text{-Score} = \frac{2 \times TP}{2 \times TP + FN + FP} \quad (7)$$

where true positive (TP) and false positive (FP) are the numbers of instances correctly and wrongly classified as positive respectively. True negative (TN) and false negative (FN) are the number of instances correctly and wrongly classified as negative respectively. Accuracy is the ratio of the number of true predictions to all instances in the set as shown

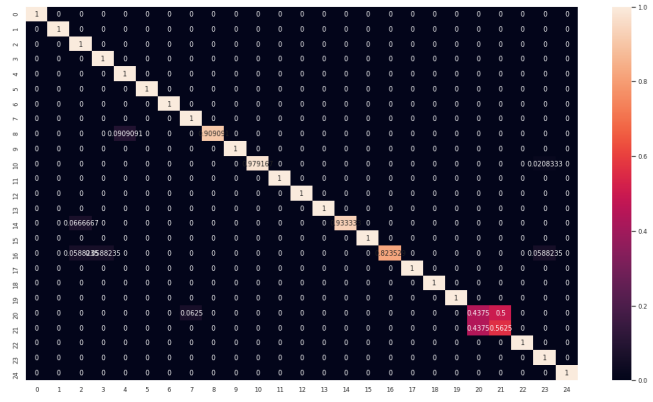
in (6). F-Score is shown as the set (7) in terms of true positives, false negatives, and false positives. Accuracy is not a correct performance metric for imbalanced datasets. On the other hand, papers compared in this work use accuracy and F-Score performance metrics to measure the success of their models. Thus, this paper gives the experiment results in terms of accuracy and the F-Score. Our main goal is that an ensemble of the CapsNet can reduce the number of FN and FP in (7).

Fig. 8 shows confusion matrices for each test part of both datasets. Each confusion matrix (Fig. 8a and 8b) implies that a model containing single CapsNet incorrectly predicts rare malware families in both datasets.

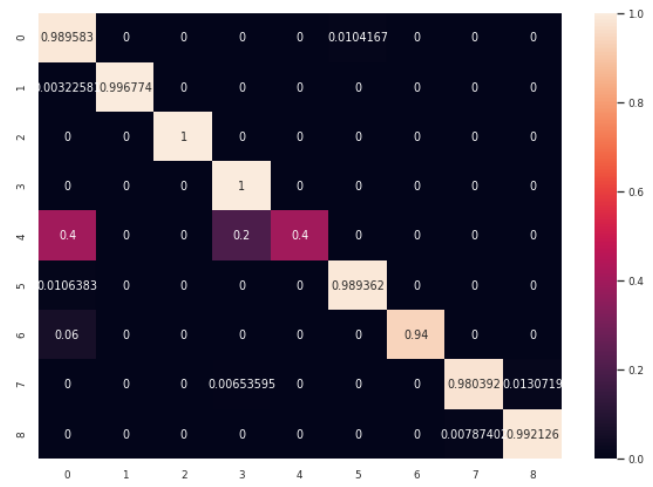
Fig. 9 is the confusion matrix of RCNF containing 5 base CapsNet models. This confusion matrix shows the prediction accuracy of the model for each malware family type in the Maling test set. Class 8, 10, 20, and 21 have been predicted wrongly by the RCNF model. On the other hand, the model has been very successful at correctly predicting other malware types in the test set. This confusion matrix also shows that RCNF is successful at correctly predicting rare malware types in the Maling test set.

In the second experiment, RCNF is tested on the BIG2015 dataset. Fig. 9 shows the prediction results of RCNF containing 10 base CapsNet for BIG2015 dataset. Class 4 is the rarest malware type in the whole dataset. Training, validation, and test sets are stratified, so the class distribution is preserved for each partition. This result shows that RCNF can predict the rarest malware type pretty well. Class 0, 1, 2, and 6 are predicted perfectly by RCNF. If the performance of RCNF is compared with the performance of a single CapsNet model, it is easily seen that RCNF is better than a single CapsNet at predicting rare malware families for imbalanced datasets.

Table 3 shows the test performance of the proposed models and others for the Maling dataset. Yue [41] uses a weighted loss function to handle the imbalanced class distribution problem in the Maling dataset; and also uses the transfer learning [61] method to classify malware family types. Due to using transfer learning, the architecture has 20M parameters; and the model is so huge. Cui et al. [42] use classical machine learning methods such as K-Nearest Neighbor and support vector machines. They have trained these algorithms using GIST and GLCM features, which are feature engineering methods for images; and they have applied resampling



(a) Maling Test Set



(b) BIG2015 Test Set

Figure 8: Confusion Matrices of single CapsNet Model for each test set.

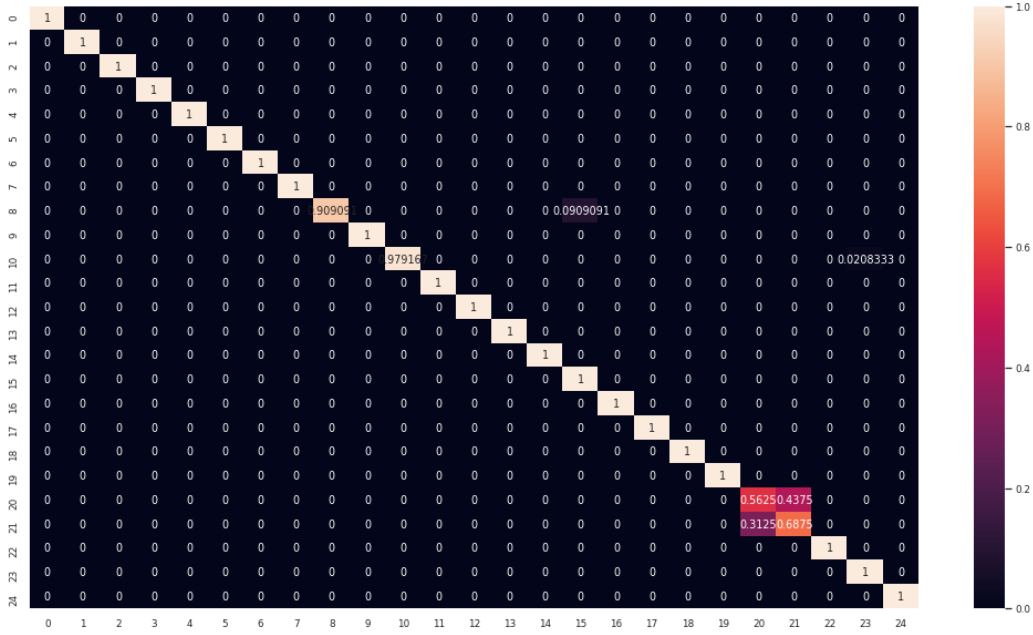


Figure 9: Confusion Matrix of 5-RCNF for Maling test set.

Table 3: Comparison RCNF and other methods for Maling test set performance.

| Model | Number of Parameters | F-Score | Accuracy |
|-----------------------------------|----------------------|---------|----------|
| Yue [41] | 20M | - | 0.9863 |
| Cui et al. [42] | - | 0.9455 | 0.9450 |
| Venkatraman et al. [43] | 212,885 | 0.916 | 0.963 |
| Vasan et al. [44] | 134M | 0.9820 | 0.9827 |
| Vasan et al. [45] | 157M | 0.9948 | 0.9950 |
| CapsNet for Maling | 90,592 | 0.9658 | 0.9863 |
| RCNF for Maling | $5 \times 90,592$ | 0.9661 | 0.9872 |

to the dataset to solve imbalanced dataset problems. RCNF does not use a weighted loss function or any sampling method to overcome the

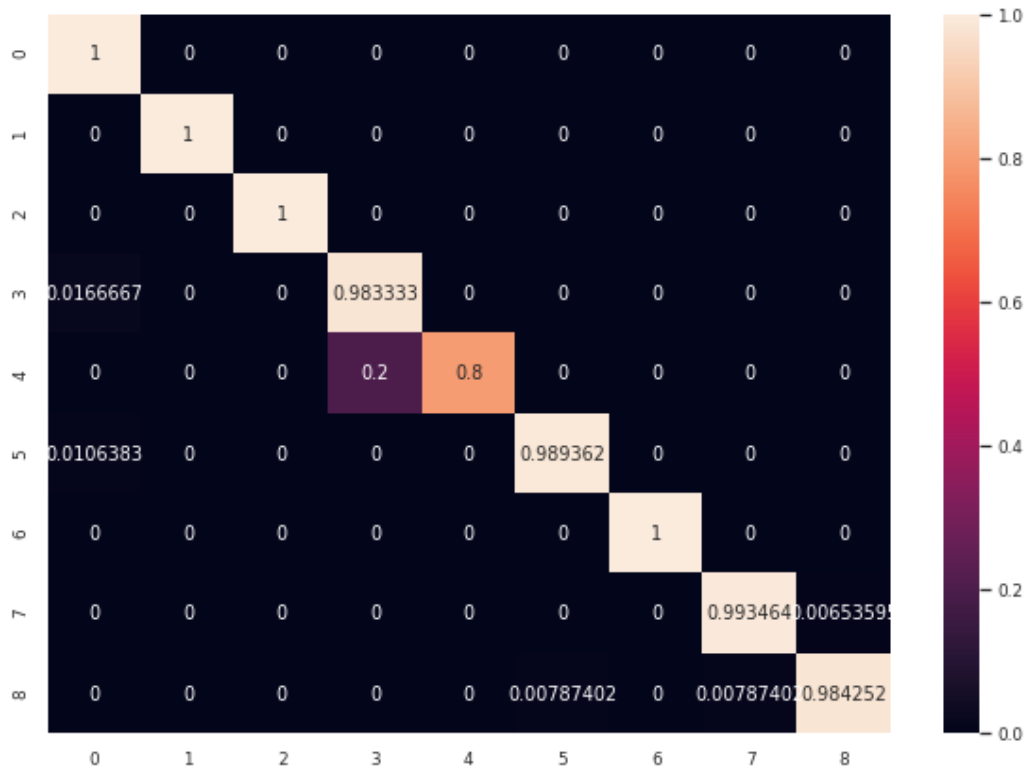


Figure 10: Confusion Matrix of 10-RCNF for BIG2015 test set.

imbalanced dataset problem. Our results are higher than these two methods; and the results also show that CapsNet and RCNF do not require any method for extra feature engineering in the Maling dataset. A single CapsNet architecture for the Maling dataset has 90,592 trainable parameters and RCNF has 452,960 trainable parameters, so our proposed methods are reasonably smaller than Yue’s method. Venkatraman et al. [43] propose two different models called CNN BiLSTM and CNN BiGRU with two variants of these models, which are called cost-sensitive and cost-insensitive. When CNN BiGRU reaches its own highest F-Score and accuracy on the Maling dataset, its number of the trainable parameters is greater than RCNF and its scores are lower than RCNF. In other words, RCNF reaches the-state-of-the-art scores in terms of F-Score and accuracy with a lower parameter size.

Table 4: Comparison RCNF and other methods for BIG2015 test set performance.

| Model | Number of Parameters | F-Score | Accuracy |
|------------------------------------|----------------------|---------|----------|
| Venkatraman et al. [43] | 212,885 | 0.725 | - |
| Cao et al. [46] | - | - | 0.95 |
| Gibert et al. [47] | - | 0.9813 | 0.9894 |
| Kreuk et al. [48] | - | - | 0.9921 |
| Le et al. [49] | 268,949 | 0.9605 | 0.9820 |
| Chen [50] | - | - | 0.9925 |
| Jung et al. [51] | 148,489 | - | 0.99 |
| Abijah et al. [6] | - | - | 0.9914 |
| Zhao et al. [52] | - | - | 0.929 |
| Khan et al. [53] | - | - | 0.8836 |
| Safa et al. [54] | - | - | 0.9931 |
| Kebede et al. [55] | - | - | 0.9915 |
| Kim et al. [56] | - | - | 0.9266 |
| Kim et al. [57] | - | 0.8936 | 0.9697 |
| Yan et al. [58] | - | - | 0.9936 |
| Naeem et al. [59] | - | 0.971 | 0.9840 |
| Jang et al. [60] | - | - | 0.9965 |
| CapsNet for BIG2015 | 527,232 | 0.9779 | 0.9926 |
| RCNF for BIG2015 | $10 \times 527,232$ | 0.9820 | 0.9956 |

IMCFN [44] is an important deep convolutional network for the Malimg dataset. When the network is trained from scratch without any data augmentation, it outperforms our RCNF implementation. IMCFN also has 126,727,705 trainable parameters. When it is compared to RCNF, IMCFN is a huge convolutional network. On the other hand, RCNF is more accurate than IMCFN, but it has a lower F-Score than IMCFN. The most important advantage of RCNF is using a lower number of parameters, and this makes RCNF trainable in GPUs with low memory.

IMCEC [45] is one of the most successful an ensemble model of CNN architectures for the MALIMG in terms of accuracy and F-Score. How-

ever, IMCEC is more complex than RCNF. IMCEC uses two different CNN architectures VGG16 [62] and ResNet-50 [63] to use transfer learning. Although IMCEC is more accurate than RCNF in terms of accuracy and F-Score, IMCEC has a total of 157M parameters because of using these CNN networks. Unlike IMCEC, RCNF has 452,960 trainable parameters and it is reasonably accurate as much as IMCEC.

Table 4 compares the test performance of the proposed models and others for the BIG2015 dataset. Venkatraman et al. [43] also test the CNN BiGRU model on the BIG2015 dataset. In this case, the performance of the CNN BiGRU model is lower than RCNF, but its number of trainable parameters are reasonably lower than RCNF. Chen [50] and Khan et al. [53] use transfer learning architectures for the dataset. Test results of our proposed methods are better than those two models, but our results are very close to Chen [50] in terms of accuracy. Our proposed models are better than Gibert et al. [47] in terms of accuracy, but the F-Score of RCNF is very close to this model. Our proposed models are better than models of Cao et al. [46], Zhao et al. [52], Kim et al. [56], and Kim et al. [57]. Jung et al. [51] propose a reasonably smaller model than our models in terms of the number of parameters, but our model has a higher accuracy score than this model. Abijah et al. [6], Safa et al. [54], Kebede et al. [55], and Yan et al. [58] propose deep learning models whose accuracy scores are close to our proposed models. In addition to these, Jang et al. [60] report five different accuracy scores for the BIG2015 dataset and they do not use F-Score despite the dataset is highly imbalanced. Their accuracy is the highest one for the dataset, but their network has two complex phases and uses GAN based data augmentation (they call this method malware obfuscator). On the other hand, RCNF does not use data augmentation, data resampling, transfer learning, and weighted loss functions for both datasets. In this perspective, RCNF is a simple version of an ensemble of CapsNet, and this simplicity highlights RCNF among its competitors.

For those tables, the last studies using Malimg and BIG2015 datasets are chosen. To compare them fairly, these models are drawn from image-based malware analysis studies. These results show that while RCNF reaches the-state-of-the-art scores in terms of F-Score and accuracy on the Malimg dataset with less trainable parameter size, it outperforms its competitors on the BIG2015 with larger size of parameters.

6. Conclusion

This paper introduces the first application of CapsNet on imbalanced malware family type classification task. Moreover, the first ensemble model of CapsNet called RCNF is introduced in this paper. The proposed models do not require any complex feature engineering methods or architecture for deep networks. To show that, we used two different malware family type datasets: Maling and BIG2015. These datasets are used for image-based malware classification. Our proposed models can utilize these datasets directly using raw pixel values.

Datasets in the paper are highly imbalanced in terms of class distribution. CapsNet and RCNF do not use oversampling, under sampling, and weighted loss function during the training phase. Results show that CapsNet and RCNF are the best models least suffering from imbalanced class distribution among others in the literature.

Experiment results show that a single CapsNet model has good performance for both BIG2015 and Maling datasets. However, we have assumed an ensemble model of CapsNet can help us to increase generalization performance and RCNF has better generalization performance results than a single CapsNet model as expected. In this point, results show that creating a bagging ensemble model CapsNet increases the performance on predicting rare malware classes. While single CapsNet can obtain 0.9779 F-Score for the BIG2015 dataset, an ensemble of 10 CapsNets achieves 0.9820 F-Score. We can observe the similar effects on the Maling dataset. It is shown that bagging increases the performance of the CapsNet for imbalanced datasets and the ensemble model is more successful at predicting rare classes than a single CapsNet model due to ensembling.

Many models compared in this paper are designed complex and large in terms of the number of parameters. Some of them use data augmentation, weighted loss functions, different extra feature engineering methods, and pre-trained deep neural networks that have large number of parameters.

As for future work, we are planning to develop a hybrid architecture for malware classification. This hybrid method will be based on CapsNet architecture.

Acknowledgement

This work is supported by The Scientific and Technological Research Council of Turkey under the grant number 118E400.

References

- [1] Najla Etaher, George RS Weir, and Mamoun Alazab. From zeus to zitmo: Trends in banking malware. In *2015 IEEE Trust-com/BigDataSE/ISPA*, volume 1, pages 1386–1391. IEEE, 2015.
- [2] Ahmad Azab, Mamoun Alazab, and Mahdi Aiash. Machine learning based botnet identification traffic. In *2016 IEEE Trust-com/BigDataSE/ISPA*, pages 1788–1794. IEEE, 2016.
- [3] Mamoun Alazab, Sitalakshmi Venkatraman, Paul Watters, and Moutaz Alazab. Information security governance: the art of detecting hidden malware. In *IT security governance innovations: theory and research*, pages 293–315. IGI Global, 2013.
- [4] Mingjian Tang, Mamoun Alazab, and Yuxiu Luo. Big data for cybersecurity: Vulnerability disclosure trends and dependencies. *IEEE Transactions on Big Data*, 2017.
- [5] Lakshmanan Nataraj, Vinod Yegneswaran, Phillip Porras, and Jian Zhang. A comparative assessment of malware classification using binary texture analysis and dynamic analysis. In *Proceedings of the 4th ACM Workshop on Security and Artificial Intelligence*, pages 21–30. ACM, 2011.
- [6] S Abijah Roseline et al. Intelligent malware detection using deep dilated residual networks for cyber security. In *Countering Cyber Attacks and Preserving the Integrity and Availability of Critical Systems*, pages 211–229. IGI Global, 2019.
- [7] Sang Ni, Quan Qian, and Rui Zhang. Malware identification using visualization images and deep learning. *Computers & Security*, 77:871–885, 2018.
- [8] Leo Breiman. Bagging predictors. *Machine learning*, 24(2):123–140, 1996.

- [9] Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton. Dynamic routing between capsules. In *Advances in neural information processing systems*, pages 3856–3866, 2017.
- [10] Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [11] Mamoun Alazab. Profiling and classifying the behavior of malicious codes. *Journal of Systems and Software*, 100:91–102, 2015.
- [12] Ahmad Azab, Robert Layton, Mamoun Alazab, and Jonathan Oliver. Mining malware to detect variants. In *2014 Fifth Cybercrime and Trustworthy Computing Conference*, pages 44–53. IEEE, 2014.
- [13] Lakshmanan Nataraj, Sreejith Karthikeyan, Gregoire Jacob, and BS Manjunath. Malware images: visualization and automatic classification. In *Proceedings of the 8th international symposium on visualization for cyber security*, page 4. ACM, 2011.
- [14] Kesav Kancherla and Srinivas Mukkamala. Image visualization based malware detection. In *2013 IEEE Symposium on Computational Intelligence in Cyber Security (CICS)*, pages 40–44. IEEE, 2013.
- [15] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [16] Daniel Nahmias, Aviad Cohen, Nir Nissim, and Yuval Elovici. Trust-sign: Trusted malware signature generation in private clouds using deep feature transfer learning. In *2019 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2019.
- [17] Amelia Jiménez-Sánchez, Shadi Albarqouni, and Diana Mateus. Capsule networks against medical imaging data challenges. In *Intravascular Imaging and Computer Assisted Stenting and Large-Scale Annotation of Biomedical Data and Expert Label Synthesis*, pages 150–160. Springer, 2018.
- [18] Parnian Afshar, Arash Mohammadi, and Konstantinos N Plataniotis. Brain tumor type classification via capsule networks. In *2018 25th IEEE International Conference on Image Processing (ICIP)*, pages 3129–3133. IEEE, 2018.

- [19] Tomas Iesmantas and Robertas Alzbutas. Convolutional capsule network for classification of breast cancer histology images. In *International Conference Image Analysis and Recognition*, pages 853–860. Springer, 2018.
- [20] Aryan Mobiny and Hien Van Nguyen. Fast capsnet for lung cancer screening. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 741–749. Springer, 2018.
- [21] Rodney LaLonde and Ulas Bagci. Capsules for object segmentation. *arXiv preprint arXiv:1804.04241*, 2018.
- [22] Ayush Jaiswal, Wael AbdAlmageed, Yue Wu, and Premkumar Natarajan. CapsuleGAN: Generative adversarial capsule network. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 0–0, 2018.
- [23] Mensah Kwabena Patrick, Adebayo Felix Adekoya, Ayidzoe Abra Mighty, and Baagyire Y Edward. Capsule networks—a survey. *Journal of King Saud University-Computer and Information Sciences*, 2019.
- [24] Solomon H Ebebuwa, Mhd Saeed Sharif, Mamoun Alazab, and Ameer Al-Nemrat. Variance ranking attributes selection techniques for binary classification problem in imbalance data. *IEEE Access*, 7:24649–24666, 2019.
- [25] Daniel Gibert, Carles Mateu, and Jordi Planes. The rise of machine learning for detection and classification of malware: Research developments, trends and challenges. *Journal of Network and Computer Applications*, page 102526, 2020.
- [26] Royi Ronen, Marian Radu, Corina Feuerstein, Elad Yom-Tov, and Mansour Ahmadi. Microsoft malware classification challenge. *arXiv preprint arXiv:1802.10135*, 2018.
- [27] Microsoft. Microsoft malware classification challenge (big 2015) — kaggle. <https://www.kaggle.com/c/malware-classification>, Feb 2015. (Accessed on 05/06/2019).

- [28] Sitalakshmi Venkatraman and Mamoun Alazab. Use of data visualisation for zero-day malware detection. *Security and Communication Networks*, 2018, 2018.
- [29] Alex Krizhevsky and Geoffrey E Hinton. Using very deep autoencoders for content-based image retrieval. In *ESANN*, volume 1, page 2, 2011.
- [30] Alberto Marchisio, Muhammad Abdullah Hanif, and Muhammad Shafique. Capsacc: An efficient hardware accelerator for capsulenets with data reuse. In *2019 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pages 964–967. IEEE, 2019.
- [31] Lorenzo Rosasco, Ernesto De Vito, Andrea Caponnetto, Michele Piana, and Alessandro Verri. Are loss functions all the same? *Neural Computation*, 16(5):1063–1076, 2004.
- [32] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [33] Leo Breiman. Bias, variance, and arcing classifiers. 1996.
- [34] Yoav Freund, Robert E Schapire, et al. Experiments with a new boosting algorithm. In *icml*, volume 96, pages 148–156. Citeseer, 1996.
- [35] J Ross Quinlan et al. Bagging, boosting, and c4. 5. In *AAAI/IAAI, Vol. 1*, pages 725–730, 1996.
- [36] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for large-scale machine learning. In *12th {USENIX} symposium on operating systems design and implementation ({OSDI} 16)*, pages 265–283, 2016.
- [37] François Chollet et al. keras, 2015.
- [38] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

- [39] Travis E Oliphant. *A guide to NumPy*, volume 1. Trelgol Publishing USA, 2006.
- [40] Wes McKinney et al. Data structures for statistical computing in python. In *Proceedings of the 9th Python in Science Conference*, volume 445, pages 51–56. Austin, TX, 2010.
- [41] Songqing Yue. Imbalanced malware images classification: a cnn based approach. *arXiv preprint arXiv:1708.08042*, 2017.
- [42] Zhihua Cui, Fei Xue, Xingjuan Cai, Yang Cao, Gai-ge Wang, and Jinjun Chen. Detection of malicious code variants based on deep learning. *IEEE Transactions on Industrial Informatics*, 14(7):3187–3196, 2018.
- [43] Sitalakshmi Venkatraman, Mamoun Alazab, and R Vinayakumar. A hybrid deep learning image-based analysis for effective malware detection. *Journal of Information Security and Applications*, 47:377–389, 2019.
- [44] Danish Vasan, Mamoun Alazab, Sobia Wassan, Hamad Naeem, Babak Safaei, and Qin Zheng. Imcfn: Image-based malware classification using fine-tuned convolutional neural network architecture. *Computer Networks*, 171:107138, 2020.
- [45] Danish Vasan, Mamoun Alazab, Sobia Wassan, Babak Safaei, and Qin Zheng. Image-based malware classification using ensemble of cnn architectures (imcec). *Computers & Security*, page 101748, 2020.
- [46] Dongzhi Cao, Xinglan Zhang, Zhenhu Ning, Jianfeng Zhao, Fei Xue, and Yongli Yang. An efficient malicious code detection system based on convolutional neural networks. In *Proceedings of the 2018 2nd International Conference on Computer Science and Artificial Intelligence*, pages 86–89. ACM, 2018.
- [47] Daniel Gibert, Carles Mateu, and Jordi Planes. An end-to-end deep learning architecture for classification of malwares binary content. In *International Conference on Artificial Neural Networks*, pages 383–391. Springer, 2018.

- [48] Felix Kreuk, Assi Barak, Shir Aviv-Reuven, Moran Baruch, Benny Pinkas, and Joseph Keshet. Deceiving end-to-end deep learning malware detectors using adversarial examples. *arXiv preprint arXiv:1802.04528*, 2018.
- [49] Quan Le, Oisín Boydell, Brian Mac Namee, and Mark Scanlon. Deep learning at the shallow end: Malware classification for non-domain experts. *Digital Investigation*, 26:S118–S126, 2018.
- [50] Li Chen. Deep transfer learning for static malware classification. *arXiv preprint arXiv:1812.07606*, 2018.
- [51] Byungho Jung, Taeguen Kim, and Eul Gyu Im. Malware classification using byte sequence information. In *Proceedings of the 2018 Conference on Research in Adaptive and Convergent Systems*, pages 143–148. ACM, 2018.
- [52] Yuntao Zhao, Chunyu Xu, Bo Bo, and Yongxin Feng. Maldeep: A deep learning classification framework against malware variants based on texture visualization. *Security and Communication Networks*, 2019, 2019.
- [53] Riaz Ullah Khan, Xiaosong Zhang, and Rajesh Kumar. Analysis of resnet and googlenet models for malware detection. *Journal of Computer Virology and Hacking Techniques*, pages 1–9, 2018.
- [54] Haidar Safa, Mohamed Nassar, and Wael Al Rahal Al Orabi. Benchmarking convolutional and recurrent neural networks for malware classification. In *2019 15th International Wireless Communications & Mobile Computing Conference (IWCMC)*, pages 561–566. IEEE, 2019.
- [55] Temesguen Messay Kebede, Ouboti Djaneye-Boundjou, Barath Narayanan Narayanan, Anca Ralescu, and David Kapp. Classification of malware programs using autoencoders based deep learning architecture and its application to the microsoft malware classification challenge (big 2015) dataset. In *2017 IEEE National Aerospace and Electronics Conference (NAECON)*, pages 70–75. IEEE, 2017.

- [56] Chang Hoon Kim, Espoir K Kabanga, and Sin-Jae Kang. Classifying malware using convolutional gated neural network. In *2018 20th International Conference on Advanced Communication Technology (ICACT)*, pages 40–44. IEEE, 2018.
- [57] Jin-Young Kim and Sung-Bae Cho. Detecting intrusive malware with a hybrid generative deep learning model. In *International Conference on Intelligent Data Engineering and Automated Learning*, pages 499–507. Springer, 2018.
- [58] Jinpei Yan, Yong Qi, and Qifan Rao. Detecting malware with an ensemble method based on deep neural network. *Security and Communication Networks*, 2018, 2018.
- [59] Hamad Naeem, Bing Guo, Muhammad Rashid Naeem, Farhan Ullah, Hamza Aldabbas, and Muhammad Sufyan Javed. Identification of malicious code variants based on image visualization. *Computers & Electrical Engineering*, 76:225–237, 2019.
- [60] Sejun Jang, Shuyu Li, and Yunsick Sung. Fasttext-based local feature visualization algorithm for merged image-based malware classification framework for cyber security and cyber defense. *Mathematics*, 8(3):460, 2020.
- [61] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? In *Advances in neural information processing systems*, pages 3320–3328, 2014.
- [62] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [63] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.