



# Deep Learning Classification for Encrypted Botnet Traffic: Optimising Model Performance and Resource Utilisation

Lucas Carr<sup>(✉)</sup>  and Josiah Chavula 

Computer Science Department, University of Cape Town, Cape Town, South Africa  
crrrluc003@myuct.ac.za, jchavula@cs.uct.ac.za

**Abstract.** Detection of malicious traffic on a network is critical to ensuring the safety and security of internet systems. Classical approaches to this task increasingly struggle with modern networking procedures, like encryption. Deep learning (DL) offers an alternative approach to traffic classification problems. We address two major problem classes: (1) botnet detection and (2) botnet family classification. For each problem, we explore five implementations of DL architectures: a multi-layer perceptron (MLP), shallow and deep convolutional neural network (CNN v1 and CNN v2), an autoencoder (AE) and an autoencoder + convolutional neural network (AE+CNN). Our evaluation of models for each respective problem class is based on the classification performance and computational requirements of each model. We further investigate the effect of training the models on an input with a reduced feature space, where we evaluate the impact this has in terms of a trade-off between computational and classification performance. For botnet detection, we find that all models attain good ( $\geq 0.979$  accuracy) classification performance on a normal testing set; however, this performance drops fairly substantially when evaluated on a set of unknown botnet families. Furthermore, we observed a clear trend between increased feature space and memory utilisation, while finding no evidence of a trend between inference time and feature space. For botnet classification, we found that models which implement CNN architectures outperform others by a substantial margin ( $\approx 6$  percentage points). We observe the same trend between feature space and memory utilisation, and absence of apparent relationship between feature space and inference time.

**Keywords:** Deep Learning · Machine Learning · Malware Classification · Malware Detection · Botnets

## 1 Introduction

The proliferation of computers and networks as tools essential to modern life has resulted in innumerable benefits. However, adoption of the associated technologies has created new security dynamics to consider; specifically in the form of

malicious software, or malware. Malware is an umbrella term that encompasses various types of software designed to infiltrate systems without permission, aiming to cause harm or exploit vulnerabilities, often with a financial motive [15].

While there are numerous sub-categories of malware, we focus attention on botnets, out of recognition that the increasing number of security-vulnerable Internet of Things (IoT) devices offer an ideal landscape for botnets [3]. A botnet defines a distributed network of computers, or *bots*, infected with software that enables the bots to be controlled by a malicious operator, or *botmaster* [1, 17]. A botnet typically leverages one of many additional types of malware - such as a worm - to propagate itself across multiple computers, and can incorporate a centralised or decentralised operating procedure [17]. Moreover, botnets attempt to hide themselves by transmitting normal traffic amongst their botnet traffic [17]. However, a defining characteristic of botnets is the presence of command and control channels, through which the malicious operator is able to transmit instructions or receive information. A common instruction would be a distributed denial-of-service (DDOS) attack, where the bots flood a target to disrupt its service [1, 3]. It is this characteristic of botnets - that the bot must at some point connect to its botmaster - that may be leveraged to build detection models. When a bot connects to the botmaster, a sequence of network flows, defined as a grouping of related traffic, can be extracted from the generated traffic, from which a deep learning (DL) model will be able to learn distinguishing patterns [17]. DL is a field within machine learning which is defined by the use of multi-layered architectures, enabling models to learn complex, hierarchical patterns from data without requiring feature engineering [12, 20].

Preventative measures against malware and botnets, are not a novel concern; there are existing approaches to detect and block malicious traffic on networks. These approaches typically deploy a Network Intrusion Detection System, or NIDS [18]. NIDS operate by implementing a broad range of techniques to detect and identify malicious traffic: notably, port analysis, blacklisted IP addresses, and inspecting packet payloads [16, 29]. Recently adopted practices around networking have dampened the effectiveness of these techniques, making NIDS which use them less reliable. Port numbers have become less reliable indicators of application type; additionally, the existence of port-obfuscation enables creators of malware to avoid detection [29]. Similarly, dynamic IP addresses and IP spoofing make systems which filter traffic based on blacklisted IPs unreliable. Finally, approaches which aim to detect malware by inspecting the payload contents of packets flowing through the network face increasing difficulty as more network traffic adopts encryption protocols - a Cisco report from 2017 noted that  $\approx 75\%$  of analysed malicious traffic made use of encryption [26].

In recognition of the shortcomings of existing malware detection practices, we define an objective to evaluate the effectiveness of different DL algorithms when tasked with detection and classification of botnet traffic using network flows. More specifically, we implement a series of binary classification models to detect malicious traffic, and a secondary series of multiclass classification models to classify botnet traffic into respective families. While this approach is not itself

novel, much of the existing literature evaluates the effectiveness of DL models in terms of the accuracy, F1-score, False Positive Rate (FPR), and False Negative Rate (FNR) [17, 19, 28].<sup>1</sup> These metrics provide insight into classification performance; however, we argue that insight into the computational requirements of a model are important. A model which attains good accuracy scores might be impractical due to its computational requirements, especially on smaller, lower-resourced networks. Moreover, models are trained on a datasets made up of only a subset of existing botnets, while newer botnets are continuously developed. Evaluating models on an unseen testing set comprised of botnet families present in their training set ignores this concept. As a result, there is greater uncertainty into a model’s ability to generalise to newer botnets.

This paper evaluates the effect of feature space size has computational performance, in terms of a model’s inference time and memory utilisation. Furthermore, a supplementary set of unknown botnet families is used to evaluate a model’s ability to generalise to zero-day attacks. More specifically, this paper makes the following contributions:

1. Evaluate the performance of five binary classification models, an MLP, shallow CNN (v1), deep CNN (v2), AE, and AE+CNN, on the standard and proto zero-day test set.
2. Implement and evaluate performance of five multiclass classification models, an MLP, shallow CNN, deep CNN, AE, and AE+CNN, which aim to identify respective families of botnets.
3. Evaluate how reducing the feature space, into 50% and 30% samples, effects the memory requirements and inference time, in relation to the overall accuracy of a model.

## 2 Related Work

### 2.1 Classification Approaches: Payload vs. Flow Based

Network traffic is made up of discrete blocks of data, called packets, which travel through a network. Approaches to classify network traffic typically make use of training data comprising of either the individual packets’ payload, or network flows. Informally, network flows represent a sequence of packets between a source and destination [29]. Statistical features can be extracted from network flows, which explain metrics such as the rate at which packets flow back and forth, and the mean packet size of the flow [7].

Another approach is to use the core contents of a packet (the payload) as training data. The notion is that the payload of malicious traffic contains at least part of the malware binary, from which a model would be able to recognise patterns belonging to this binary [11]. Payload based approaches face the difficulty of classifying encrypted traffic - a problem that flow-based analysis avoids, since only the packet headers are required, which are not encrypted [12]. There

---

<sup>1</sup> Definitions for these are found in Sect. 5.1.

have been implementations of payload-based classifiers which are able to handle encrypted traffic [4,10,11]. These approaches typically require thorough processing steps to prepare the data for classification, which makes payload-based approaches ill-suited to real world application. Conversely, aggregated network flows are comparatively easy to extract [18].

## 2.2 Machine Learning Approaches

Hadidi et al. [7] evaluate the effectiveness of different machine learning approaches to botnet detection. Their approaches include Support-Vector Machines (SVM), K-Nearest Neighbour (KNN), and Bayesian Networks (BN) to classifying botnet traffic based of either payloads or network flows, using Detection Rate (DR) and False-Positive Rate (FPR) as evaluation metrics.<sup>23</sup> Simulated network traffic was captured in a sandbox environment. Non-encrypted traffic was used, so as to make their payload-classifier able to handle the traffic captures. Notably, in the network flow preprocessing phases, identifying features such as IP addresses and port numbers were removed from the datasets [7]. In the majority of evaluations, payload-based classifiers have been shown to be superior to flow-based methods. Specifically, the payload-based models such as KNN, SVM, and BN have recorded detection rates (DRs) of 1, 0.995, and 0.938 respectively. In contrast, the flow-based models have posted comparatively lower scores, with DRs of 0.968, 0.910, and 0.838. This trend of better performance by payload-based classifiers is also evident in terms of false positive rates (FPRs).

Yeo et al. [28] evaluate four different ML architectures (Random Forest, CNN, MLP, and SVM) to be used as binary classifiers for botnet detection. The models were trained on bi-directional network flows extracted from PCAPs in the CTU-13 dataset - a dataset containing botnet traffic from 7 different families. Typical measurements of accuracy, precision and recall were used as evaluation metrics, while the performance of a classifier was evaluated w.r.t an individual botnet family.

## 2.3 Autoencoders

Autoencoders (AE) are a type of unsupervised learning algorithm that aim to compress input data into a lower-dimensional “latent” vector. As output, an approximation of the input is reconstructed from the latent vector, through a decoding process [6]. When applying AEs to classification tasks, the decoding process can be replaced with a classification layer (like a softmax or sigmoid layer). The encoding portion of the AE focuses on detecting crucial features in the input data and encoding them into a condensed representation. This process of dimensionality reduction ideally encodes the most significant features which aid the subsequent classification layer.

---

<sup>2</sup> Detection Rate is identical to Recall.

<sup>3</sup> DR and FPR defined as  $\frac{TP}{TP+FN}$ ,  $\frac{FP}{FP+TN}$ , respectively.

Deep Packet, an approach proposed by Lotfollahi et al. [11], is a system which incorporates both feature extraction and classification stages. This approach is not directly related to malware detection or classification, and instead aims to identify major traffic classes or application. While this is a significant divergence from the aims of our paper, the approach to solving their problem using DL has strong parallels to ours. They propose a five-layered Stacked Auto Encoder (SAE) connected to a softmax layer, and 1D-CNN as classifiers made up of two convolutional layers and a softmax layer [11]. Following convention, Precision, Recall and F1-score were the chosen evaluation metrics.

## 2.4 Convolutional Neural Networks

Convolutional neural networks (CNNs) have become increasingly popular algorithms for classification [14]. Typically CNNs work with grid-like inputs, such as an image, where a convolutional operation will sweep over the grid, producing a feature map which represents significant areas of the input. These feature maps enable the model to learn identifying spatial patterns in data. Most applications of CNNs use two-dimensional image data, or image-representations of streams of one-dimensional data. However, the fundamental principle of a sequence of convolutional layers that identify increasingly complex patterns in the data holds for inputs which are not grid-like in nature - for instance, a network flow, which is a one-dimensional vector [2].

Marín et al. [12] used flow-based and packet-based approaches to detect (a binary classification) and further classify (a multiclass classification) botnet traffic. These approaches implement a 1D CNN which is connected to an LSTM. For the binary classification task, the flow-based approach is the best model by a significant margin, with an accuracy of 0.986, compared to the packet-based model's 0.776. They note the flow-based model is able to achieve this accuracy with a FPR of  $\approx 0.025$ . For the multiclass classification, they were unable to use a flow-based approach due to limitations relating to their dataset. Bearing this in mind, their packet-based model, which aimed to classify traffic into classes of Benign, Neris, Rbot, and Virut, attained accuracies of 0.878, 0.635, 0.999, and 0.547 for each respective class, with an overall accuracy of 0.765.

Pektas & Acarman [17] proposed using a deep neural network (DNN) as a binary classifier for botnet detection. They employed the CTU-13 dataset for botnet captures, which was also used in Yeo et al. [28]. To process the data, they constructed a graph representation of the network captures where nodes represent connected hosts. This approach allowed them to extract statistical information about network flows. Alongside source and destination IP addresses and port numbers, they computed five statistical metrics for each flow: mean, median, maximum, minimum, standard deviation. These metrics were applied to the duration, byte size, number of packets and periodicity of each flow. For evaluation, they focused on accuracy, precision, recall, and the F1-score.

### 3 Datasets and Preprocessing

Detection and classification of botnet traffic using DL algorithms is a task that lends itself towards supervised learning. Supervised learning requires the use of high quality, labelled datasets with sufficient samples for the training and evaluation process - the existence of such datasets is rare [26].

Network traffic is typically captured through programs like WireShark, where the information is stored in PCAP files. For the task at hand, network flows, and ideally bidirectional network flows, are extracted from these PCAP files, and preprocessed into suitable training data. We developed a preprocessing pipeline which received traffic captures in the form of PCAP files as input, and after a series of steps, outputted datasets in the form of .csv files. Information concerning the original source of the dataset is discussed in Sect. 3.1, after which Sect. 3.2 describes the preprocessing steps taken in the pipeline.

#### 3.1 Dataset

The Stratosphere Research Laboratory host an online repository<sup>4</sup> of malicious and normal network captures. Specifically, they have created the CTU-13 dataset, which contains network captures of real traffic from seven distinct botnet families [5]. The dataset is made up of thirteen captures; each capture containing the malware binary, extracted network flows, and a PCAP file containing only botnet traffic from that scenario’s capture (these PCAP files have had their normal and background traffic removed due to privacy considerations) [21].

The bidirectional network flows provided by the CTU-13 dataset, which were extracted using the open source tool, openArgus, offer comparatively limited information, relative to what could be extracted when using CICFlowMeter [8].<sup>5</sup> Consequently, only the PCAP files containing botnet traffic were used from this dataset. These were then supplemented with the Stratosphere Research Laboratory’s repository of normal captures, which are captures of network activity which imitate a typical user’s activity on a network, and are restricted to contain only benign network activity. The inclusion of benign traffic was necessary in order to facilitate the measurement of True Negatives and False Positives; as well as encourage models to be able to generalise to a real-world environment [22].

Captures 10 and 11 were omitted from the CTU-13 collection, as they were instances of a malware family which had sufficient representation from the remaining scenarios. The resultant eleven scenarios were supplemented with five ‘normal’ captures.

For the binary classification process, the botnet families Murlo and NSIS-ay were excluded from the training and testing sets. This was to enable the creation of an additional testing set, hereafter referred to as the ‘proto zero-day’

<sup>4</sup> The repositories can be found at: [www.stratosphereips.org/datasets-overview](http://www.stratosphereips.org/datasets-overview).

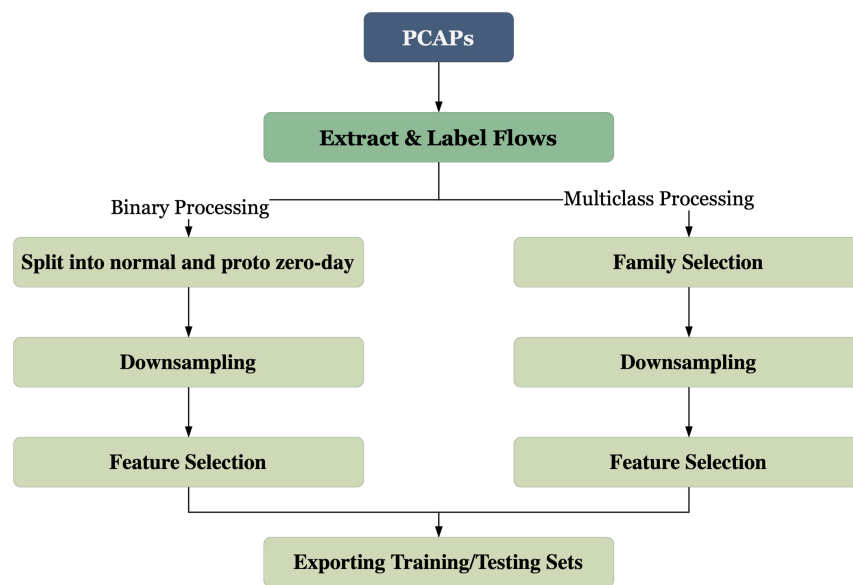
<sup>5</sup> openArgus can be found at: <https://openargus.org>.



set, which contained botnet families to which the model had not been exposed. Unlike traditional test sets, which present models with unseen instances of known botnet families, our set introduces entirely new categories. This additional measure is analogous to concept of zero-day attacks, which are malware attacks that have never appeared before [30]. While measuring a model’s exact ability to detect zero-day attacks would be impossible, we argue that this approach reasonable indication of the model’s performance when encountering previously unseen attacks.

### 3.2 The Pipeline

Illustrated in Fig. 1, the preprocessing pipeline begins with a collection of PCAP files representing traffic captures. These were the eleven botnet captures and five normal captures. Individually, each PCAP represents a capture of either entirely normal traffic, or a single botnet family [5]. The bi-directional network flows were extracted using CICFlowMeter [8]. The extracted flows from each PCAP file would be stored in a corresponding .csv file.



**Fig. 1.** Illustration of the preprocessing pipeline, showing the divergence in processing steps for the binary and multiclass classification datasets

**Flow Extraction and Labelling.** An essential part of the preprocessing was to assign accurate labels to the data. The CTU-13 dataset includes labelled bi-directional network flows for each scenario, extracted with the openArgus. However, these flows are not as detailed as flows extracted using an alternative tool: CICFlowMeter [8]. A consequence is that flows extracted using CICFlowMeter needed to be labelled manually. The nature of the sourced PCAP files was that they contained either entirely malicious or entirely benign traffic [21]. As such, the labelling process was straightforward to implement: the labels corresponding

to the flows generated from the previous stage could be identified by knowing which PCAP file the flows originated from - which was simple to do, given that each PCAP file would produce a single .csv of its extracted flows. The labelling process was automated through a python script, *FlowLabeller.py*.

For the binary classification dataset, the data was either labelled as 0 indicating benign, or 1 meaning malicious. For the multi-class classification, benign traffic was labelled 0, and the malware classes were labelled from 1-7.

**Feature Selection.** The bidirectional flows are one-dimensional vectors, made up of 82 features. Each feature is a specific measurement of how the data behaves in the flow, from which patterns can be learned during the training process. For example, there is a feature (Total Fwd Pkt) which provides the total number of forward flowing packets in the bi-directional flow. However, we recognised that allowing certain features to persist in the dataset could potentially be detrimental to the models' classification performance.

Features relating to IP addresses and port numbers were removed. Dynamic IP address, port-obfuscation and IP spoofing are techniques which make relying on these features for classification a poor idea [7, 18]. Additionally, the models should be able to generalise to unseen data as best as possible, and inclusion of these features in the training set is antithetical to this goal, since they are not intrinsic to the identity of the malicious traffic [18].

In the end, each bi-directional flow is represented as a one-dimensional vector with 75 features. We then create two additional datasets containing a random sample of 50% and 30% of the features (37 and 22 features, respectively). This would facilitate investigation into how reducing the feature space might lower memory requirements and inference time, and what effect it would have on classification performance.<sup>6</sup>

**Balancing.** The datasets for the binary classification task were balanced to have an even distribution of benign and malicious samples. The malicious samples were made up of Neris, Rbot, Virut, Menti, and Sogou botnet families. Murlo and NSIS.ay families were excluded as they were used for the creation of the proto zero-day dataset. The approximately 110,000 malicious flows were down-sampled to 59,000, which was the number of benign samples in the dataset; the datasets were split into training, validation, and testing sets in a 72%, 8%, and 20% ratio.

To balance the classes in the datasets for multiclass classification, we ensured that classes would have a sufficient number of respective samples, so as to allow the model to train well on that class. Empirically, we determined that a class required a minimum of 30,000 samples in order for the classifiers to perform effectively. A consequence of this being that botnet families Sogou, Menti, NSIS.ay, and Murlo were removed from the training set. As indicated in Table 1, these families did not have enough samples for practical up-sampling. The resultant

---

<sup>6</sup> The specific features present in these datasets are described in Table 2 in the Appendix.



**Table 1.** Botnet Families and Network Flows from the CTU-13 Dataset with additional Benign Traffic

Family	#Flows	Family	#Flows
Neris	190,028	Sogou	72
Rbot	46,796	dMurlo	11,537
Virut	85,779	NSIS.ay	7,645
Menti	4,810	Benign	56,665

dataset included traffic labelled as benign, Rbot, Virut, or Neris. These classes were then down-sampled to consist of 45000 samples each; the datasets were split into training, validation, and testing sets in a 72%, 8%, and 20% ratio.

## 4 Implementation

### 4.1 Hyperparameter Tuning

The performance of a DL model is heavily influenced by hyperparameters. Discovering optimal hyperparameters typically involves references to existing literature, and exploring iterations of training slightly different models and evaluating which parameters yield better results (for example, Grid-Search). This process is both computationally expensive and time consuming. We adopt an alternate approach using an extension of Keras Tuner called Hyperband [9,13]. Starting with a predefined set of options, including ranges of layer sizes, activation functions, learning rates, and dropout rate, Hyperband adopts an early-stopping strategy to identify promising combinations of hyperparameters. These are trained for a small number epochs to assess their performance. The top-performing configurations are kept for further training, while the rest are discarded. This process is iterated until the algorithm converges on a network topology and set of hyperparameters that yield near-optimal performance [9].

### 4.2 Architectures

To achieve our aims of both detecting and classifying botnet traffic, we decided that for each model, we train a binary classifier (for botnet detection) and a multiclass classifier (for botnet classification). This section describes the topology and hyperparameters of each model - arrived at through implementation of the Hyperband process discussed in Sect. 4.1.

**Multilayer Perceptron.** The MLP is, by design, our simplest model. The MLP binary classifier has an input layer, connected to a single densely connected layer with 33 neurons using the Tanh activation function. The output of this layer is fed into another densely connected layer with a Sigmoid activation function for

classification. The architecture for the multiclass classification model is markedly similar, the difference being an increase in the size of the hidden layer, with 128 neurons, and a classification layer which uses a Softmax function.

**Convolutional Neural Networks.** For both classification tasks we introduced two CNN architectures inspired by the implementations of 1D-CNNs as per [25, 29]. Each task has a respective shallow CNN (CNN v1), and deeper CNN (CNN v2). Across all models, we adopted the Adam optimiser during the training phase, which has had widespread success in related literature [11, 26, 29]. Furthermore, all networks shared a common filter size of  $3 \times 1$ , with a stride of 1. Following these convolutional layers, MaxPooling was employed as a down-sampling technique to reduce spatial dimensions and retain critical features of the input.

With respect to the binary classifiers, CNN v1 had two 1D convolutional layers made up of 128 and 416 filters, respectively. The output from the final MaxPooling layer was flattened, and fed into a densely connected layer with 352 neurons. A final Sigmoid layer was used for classification. CNN v2 implemented three 1D convolutional layers, with 40, 136, and 232 filters. After the final MaxPooling layer, a dropout of 0.5 was introduced to combat overfitting. The result was flattened, and channelled into a dense layer of 104 neurons, followed by another dense layer of 40 neurons before a final Sigmoid layer for classification. Aside from the classification layer, all applicable layers made use of the Rectified Linear Unit (ReLU) activation function, which introduced non-linearity to the model - a decision determined through the hyperparameter tuning process.

For multiclass classifiers, CNN v1 had two 1D convolutional layers made up of 232 filters each; after the second MaxPooling layer, the output was flattened and inputted to densely connected layer of 72 neurons, before a final Softmax classification layer. CNN v2 implemented three 1D convolutional layers, made up of 232, 104, and 40 filters, respectively. After the third MaxPooling layer, we introduced a 0.5 dropout to the model for overfitting. The output was then flattened and inputted to a densely connected layer of 296 neurons, after which another dropout layer of 0.25 was introduced. Two more densely connected layers sized 456, and 168 were implemented before the Softmax classification layer. The dense and convolutional layers used the Tanh activation function - the decision to implement Tanh was made through empirical findings, through the hyperparameter tuning process.

**Autoencoders.** The architecture of the AEs we implemented for the binary classification and multiclass classification problems were very similar, with differences appearing in the size of the layers. Drawing inspiration from [11], each AE had five fully connected hidden layers, and a classification layer (Sigmoid or Softmax). The sizes of these layers for the binary classifier were [232, 72, 40, 104, 232], whereas the multiclass classifier has layers sized [168, 104, 104, 40, 104]. Both the binary and multiclass classifiers implement a Tanh activation function in each layer.

**AE+CNN.** The AE + CNN is an ensemble of the previously implemented AE and CNN v1. For the binary classification task, five densely connected layers were used for the encoding process, which had 136, 136, 72, 104, 136 neurons, respectively. The output from the fifth layer was reshaped in order to be suitable input for the CNN. Subsequently, two one-dimensional convolutional layers were implemented with 64 and 32 filters, respectively. Each convolutional layer was followed by a MaxPooling layer, where the final MaxPooling layer was flattened and channelled into a densely connected layer with 8 neurons, connected to the final classification layer which used the Sigmoid activation function. All of the applicable layers used a ReLU activation function, and the Adam optimiser - for the same reasons as before.

The multiclass AE+CNN implemented a similar architecture, with five densely connected encoding layers with 136, 104, 72, 40, 104 neurons, respectively. The same reshaping and subsequent convolutional and MaxPooling layers were included, with 32 and 96 neurons in each respective convolutional layer. These layers, where applicable, implemented a Tanh activation function, as opposed to the binary classification model’s ReLU. A final softmax layer was used for classification.

## 5 Experiment Design

### 5.1 Evaluation Metrics

To evaluate the performance of a model, we use accuracy, FPR, FNR, mean inference time (MIT), and mean memory usage (MMU). While accuracy is a standard metric w.r.t. determining the effectiveness of a classifier, much of the related work prefers F1-score [10, 11, 18]. The advantage F1-score offers is that it provides a fairer representation of a model’s performance in the case of unbalanced datasets [27]. In our case, measures were undertaken to ensure a balanced training and testing set; consequently, accuracy was preferred to F1-score.

FPR quantifies the fraction of benign results incorrectly identified as malicious by the model [7]. Conversely, FNR measures the proportion of actual threats misclassified as benign. In the context of intrusion detection systems, both metrics are important. High FNRs undermine the essential purpose of an IDS - to detect threats. On the other hand, an elevated FPR can erode trust in the system. If users are frequently alerted to false threats, they may begin to ignore genuine threats [7]. These metrics are defined as:

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$

$$\text{FPR} = \frac{FP}{FP + TN}$$

$$\text{FNR} = \frac{FN}{FP + FN}$$

where  $TP$  refers to true positives,  $TN$  to true negatives,  $FP$  to false positives, and  $FN$  to false negatives.

MIT and MMU provide the mean time for a model to make an inference, and memory required to make 100 inferences, respectively. These measurements require multiple iterations of measurements for each model, in order to extract the mean. An alternate approach of taking the worst-case performance of these measurements was considered. The advantage of a worst-case measurement is that it enables us to infer the hardware specifications needed to implement the model without fear of failure, providing an upper bound on the memory usage or inference time. However, we recognise that there are significant difficulties in ascertaining accurate measurements of memory usage or CPU time, (see Sect. 5.2 for further discussion of these challenges).

## 5.2 Binary Classification Task

### Classification Performance on Normal vs Proto Zero-Day Test Sets.

This experiment establishes a baseline evaluation of how MLP, CNN v1, CNN v2, AE, and AE+CNN perform, in terms of classification accuracy, FPR, and FNR on the conventional testing set. Subsequently, these models are evaluated on the additional testing set, as discussed in Sects. 3.1 and 3.2 to ascertain their ability to generalise to unseen botnet families. Accuracy is determined through the use of the Keras framework’s ‘evaluate’ function. For calculation of FPR and FNR, a model makes predictions on a testing set which are compared with the set of ground truths to determine the FP, FN, TP, TN values used in the formulas outlined in Sect. 5.1.

**Effect of a Reduced Input Feature Space on Computational and Classification Performance.** While DL has largely alleviated the necessity for feature engineering, as is present in machine learning, larger feature spaces typically incur a greater computational cost [11, 20]. We aim to explore the relationship between computational and classification performance of the five DL models when trained on inputs of 100%, 50%, and 30% of the feature space.

We define computational performance as the memory usage (MMU) and inference time (MIT) of a model. Determining accurate values for these metrics presents significant difficulties: during the execution of a program, there are invariably other processes running concurrently. Furthermore, memory management of an operating system is largely beyond our control. Empirically, we found that when a program iterated over each model, measuring memory use and inference time, there was a consistent increase in memory consumption with each subsequent iteration, irrespective of model complexity. To minimise the potential impact these factors may introduce, we elected to measure the inference time and memory utilisation across batches ( $b = 10$ ) of the test dataset, each batch composed with fixed number of samples ( $n = 100$ ). In this approach, the system was rebooted after each subsequent evaluation of a batch. Memory use was measured

using the ‘Memory Profiler’<sup>7</sup> python package, which provides a list of memory usage taken at snapshots during the program’s execution. We record the model and memory usage for each batch, extracting the mean usage from these results. Inference time was determined using the ‘timeit’ package from Python’s Standard Template Library [23]. We measure and record the time taken for a model to make predictions over a batch, extracting the mean from these records. We implement this procedure three times for each model, using training and testing sets containing 100%, 50%, and 30% of the available features.

### 5.3 Multiclass Classification Task

**Multiclass Classification Accuracy for Each Botnet Family.** This experiment evaluates the MLP, CNN v1, CNN v2, AE, and AE+CNN models to determine their overall accuracy scores, as well as their accuracy scores for each specific botnet family. Metrics like FPRs and FNRs are not applicable to multiclass classification problems, as they require a binary relation. It is still useful, however, to determine how each model performs relative to each class in the multiclass classification. We evaluate each model on the normal testing set only, and from this evaluation we are able to determine a general accuracy for each model, as well as the accuracy of each model respective to the available classes. Reiterating the discussion from Sect. 5.1, we use accuracy because we have ensured that each class has an equal representation of samples, at 8000. For each model, we make classifications using the Keras Library’s ‘predict()’ method, and store the results in a confusion matrix.

**Evaluating the Effect of a Reduced Input Feature Space.** We evaluate how datasets with reduced feature spaces influence the computational and classification performance of models. More specifically, five models are trained on three datasets containing a random sample of 100%, 50%, and 30% of the total features - this translates to datasets with 74, 37, and 22 features, respectively. We then observe the effect that a reduced feature space might have on a single, or group of models’ computational and classification performance. We regard the computational performance to be the MIT and MMU of a model, while classification performance is largely defined as a model’s accuracy across all classes, with a secondary focus on the portion incorrectly classified traffic. This latter focus enables us to determine which models might struggle to classify specific families, or if certain families are poorly classified by all models.

We determine the MIT and MMU by taking measurements over a series batches ( $b = 10$ ) containing a fixed number of samples ( $n = 100$ ).

---

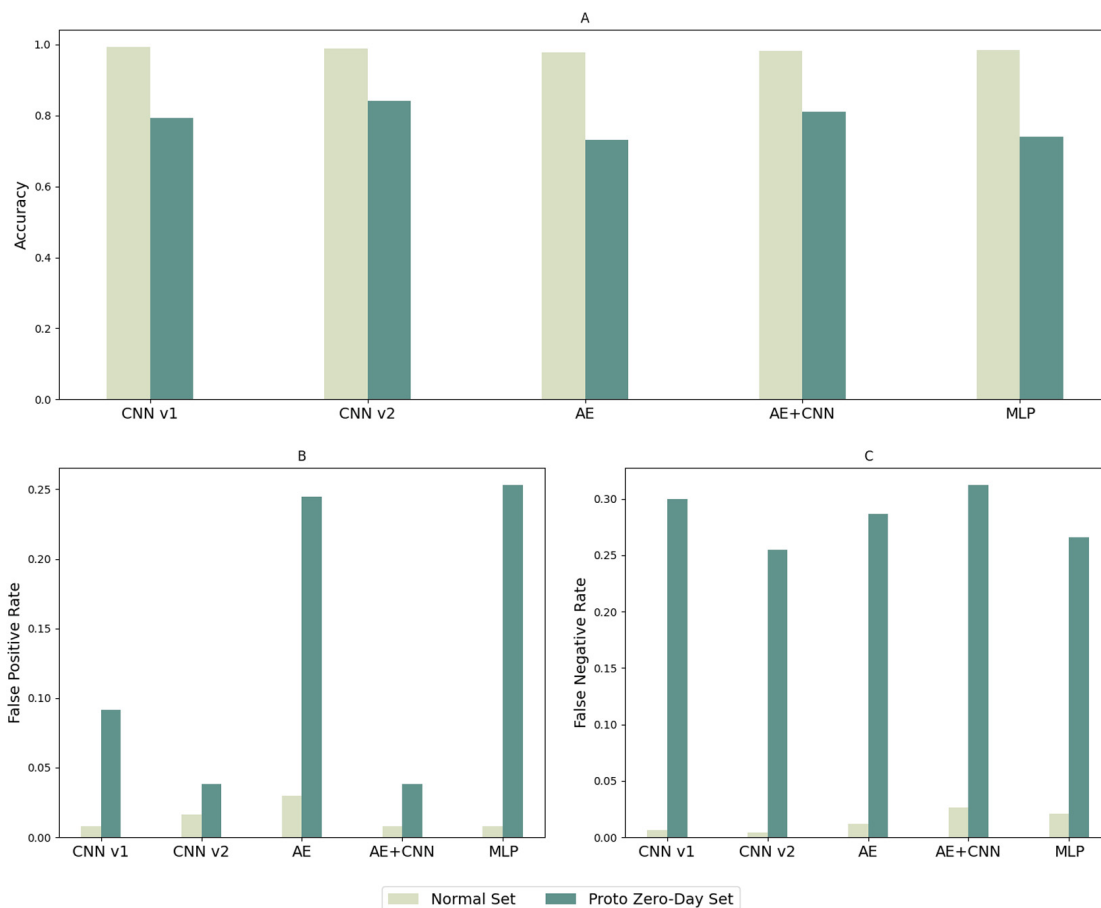
<sup>7</sup> Memory Profiler can be accessed at [https://github.com/pythonprofilers/memory\\_profiler](https://github.com/pythonprofilers/memory_profiler).

## 6 Results and Discussion

In the following, we present and discuss the findings from the experiments outlined in 5. We begin with the discussion of binary classifiers, followed by multi-class classification results.

### 6.1 Binary Classifiers for Botnet Detection

**Performance on Normal and Proto Zero-Day Test Sets.** Figure 2 displays the accuracy, FPR, and FNR of the MLP, CNN v1, CNN v2, AE, and AE+CNN when evaluated on the normal and proto zero-day testing sets. Performance on the normal testing set provides an indication of a model’s ability to generalise to unseen traffic that belongs to botnet families which were present in the training set. On the other hand, results relating to the proto zero-day testing set relates to a model’s ability to classify traffic from botnet families which were entirely excluded from the training set.



**Fig. 2.** Sub-Figures A, B, and C, respectively, show Accuracy, False-Positive Rates, and False-Negative Rates of each model when evaluated on the normal and proto zero-day testing sets.



All models achieved relatively high ( $\geq 0.979$ ) accuracy scores when evaluated on the normal testing set. The MLP, our least complex model, attained a score of 0.986, indicating that the task of botnet detection from known families is fairly simple. The CNN v2 attained the highest classification accuracy of 0.993, while the AE was the lowest with a score of 0.979. A possible explanation for the poor performance of the autoencoder models is that latent vectors created by the encoding phase fail to capture some distinguishing features of the input data. When comparing best and worst models - CNN v2 and AE - there is an absolute improvement by the CNN v2 of 0.014. While this improvement is small, we calculate that the error rate of the CNN v2, at 0.007, represents a 0.667 reduction in errors relative to the error rate of the AE, at 0.021.<sup>8</sup>

To recognise the significance of this reduction, we need to contextualise the problem. Networks are often required to handle a large volume of traffic; for example, a network on a college campus with 10,000 users may see a typical transfer of 7TB of data every 24 hour period [24]. For an intrusion detection system monitoring this network, a relative reduction in error rate of 0.667 could entail thousands of fewer errors.

These errors can be broken down into two categories: false negatives and false positives. From this, we determine the FPR and FNR for each model. When FPR was evaluated on the normal testing set, the CNN v2, AE+CNN, and MLP were the best performers with FPRs of 0.008. The CNN v2 attained an FPR of 0.016. The poorer performance of the CNN v2 in comparison to the CNN v1 might be explained by the increased complexity of CNN v2, with the additional convolutional layer and fully connected layers causing overfitting. The worst performer w.r.t FPR was once again the AE, with a measure of 0.030. It is difficult to suggest an acceptable tolerance for FPR and FNR: factors around the type of information being secured, the size of the organisation, and general security posture may all influence what might be deemed acceptable. However, the FPRs achieved by the CNN v2, AE+CNN, and MLP, with respective accuracies of 0.993, 0.983, and 0.986, present an improvement on existing work [12, 28].

Evaluation on the proto zero-day testing set showed a decline in performance in terms of accuracy, FPR, and FNR across all models. The best performing model w.r.t. this classification performance was the CNN v2, which achieved an accuracy of 0.842, FPR of 0.038, and FNR of 0.255. The same model evaluated on the normal testing set achieved scores of 0.990, 0.016, and 0.004 for each respective metric. The decline in classification performance aligns with our understanding that distinct botnet families are likely to exhibit, at least partially, different behaviour. It is this difference which causes the models to struggle when classifying traffic belonging to families entirely excluded from the training set. However, the model's performance on the proto zero-day testing set remains a significant improvement on guessing. We suggest that an explanation for this improvement rests on the notion that while there are certainly some differences, distinct botnet families must express certain shared behaviour that is not present in benign traffic. This would be behaviour intrinsic to all sampled

---

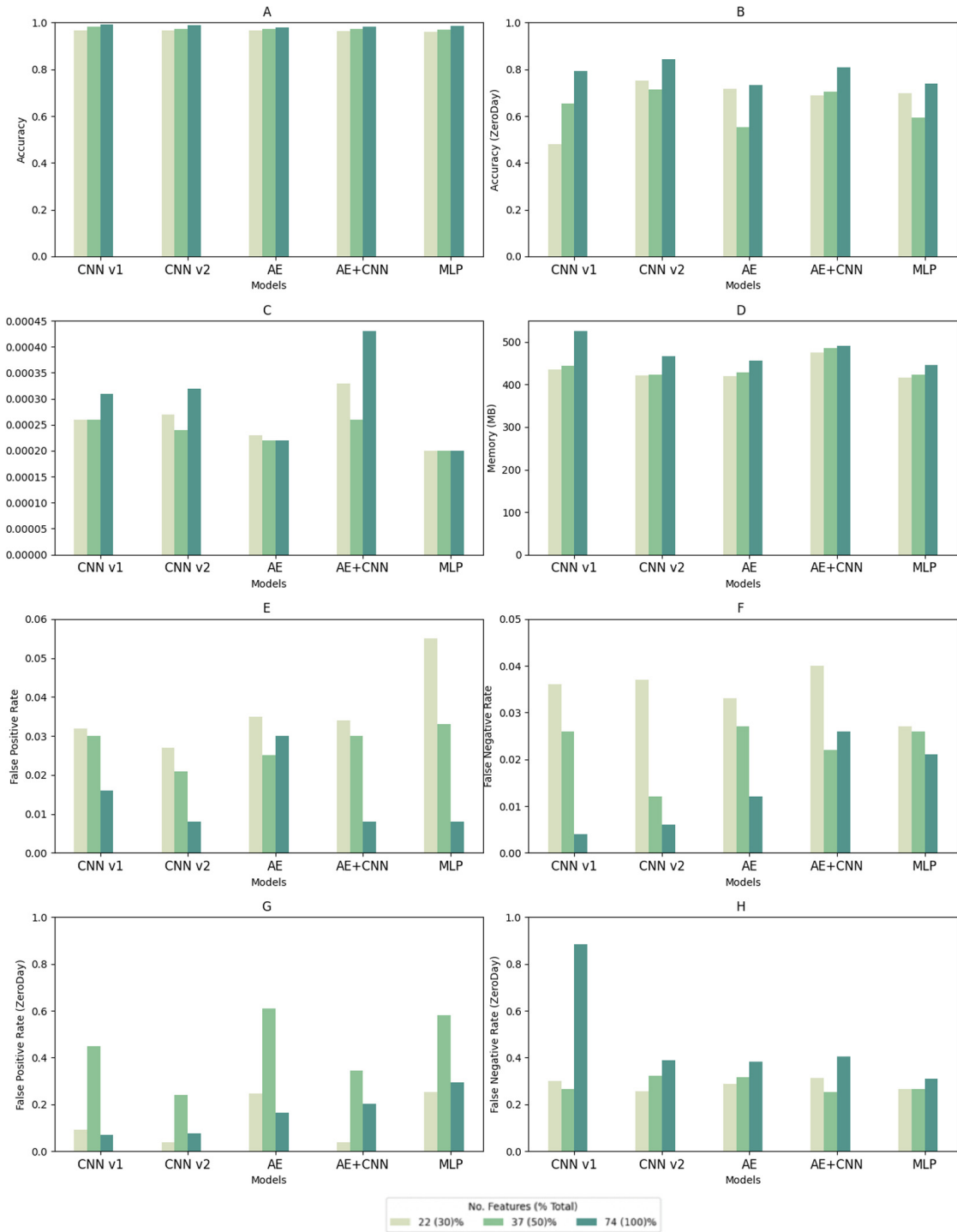
<sup>8</sup> Formulas and calculations provided in Appendix.

botnet families, as opposed to behaviour distinct to individual families. This shared behaviour among botnet families is what a model would recognise in the proto zero-day testing set.

There was also a more pronounced degree of variability in classification performance of the five models when evaluated on the proto zero-day testing set, relative to the evaluation based on the normal testing set. For instance, the mean accuracy, over all five models, when evaluated on the proto zero-day testing set was 0.783 with a standard deviation of 0.047. In contrast, when evaluated on the normal testing set, we determined a standard deviation of 0.005 around a mean of 0.986. This variability indicates that there is a fairly substantial difference between each models' ability to learn the more complex, intrinsic botnet behaviour which enable better detection of unknown botnet families.

From these two ideas - that good performance on the proto zero-day testing set requires learning more complex behaviour patterns intrinsic to all botnets, and that there is greater variability in model's classification performance when evaluated on the proto zero-day testing set - we make the claim that certain models, specifically models which implement convolutional layers, are significantly more capable of learning these more complex patterns. We ground this claim in the observation that the CNN v1, CNN v2, and AE+CNN achieved accuracies on the proto zero-day set of 0.842, 0.793, and 0.810, respectively, whereas, the AE and MLP achieved accuracies of 0.732 and 0.740. A potential explanation for the efficacy of convolutional layers towards learning these behaviour patterns is that they are particularly good at learning patterns which emerge from the relationship between closely related features - which may represent the more intrinsic behaviour general to all botnets. The AE and MLP are able to achieve high ( $>0.979$ ) accuracies on the known botnet families because they can learn the defining features of each individual family, as opposed to learning some underlying pattern seen in all families. This is sufficient for binary classification on known families, but generalises poorly to classification of unknown families, because these 'defining' features may not be present.

**Effect of a Reduced Input Feature Space on Computational and Classification Performance.** Sub-Fig. 3C and 3D show the MIT and MMU for each model when using 30%, 50%, and 100% of the available input feature space. From Fig. 3D, we observe a small but clear trend whereby increasing the feature space of a model's input has an associated increase in the memory usage of that model. This result aligns with the position outlined by Sarker [20], that the absence of feature engineering in DL may increase computational requirements of DL models. While the trend is consistent across all models, the proportional increase in memory is small; going from feature spaces of 30% to 100% (that is, 22 to 74 features), we observe a mean increase in MMU of  $\approx 10\%$ . In the most extreme case, CNN v1, the jump from 30% to 100% of features saw an increase in MMU of  $\approx 90$  MB, or 20.726%. For CNNs, an explanation for the relatively small increase to MMU when given larger feature spaces is likely found in their use of sparsely connected layers and parameter sharing, which reduce the number of trainable parameters in a model.



**Fig. 3.** Figures illustrating model performance on testing sets comprised of 30% (22 features), 50% (37 features), 100% (74 features) of total features. The figures show Normal Accuracy (A), Proto Zero-Day Accuracy (B), Mean Inference Time (C), Mean Memory Usage (D), False Positive Rate (E), False Negative Rate (F), False Positive Rate on Zero-Day set (G), and False Negative Rate on Zero-Day set (H).

From sub-Fig. 3D, we observe that the CNN v1 sees the largest increase in memory utilisation when evaluated on larger feature spaces. However, larger feature spaces in this model also result an improvement to classification performance. From sub-Fig. 3A, the normal accuracy improves from 0.968 to 0.993, sub-Fig. 3E shows that the FPR improves from a rate of 0.027 to 0.008, and FNR (sub-Fig. 3F) improves from a rate of 0.037 to 0.006. These improvements to classification performance when using larger feature spaces are significant. Moreover, they are, to lesser extents, observable in the four other models. Consequently, we make the claim that the increase in memory utilisation associated with larger feature spaces is justified by the improvements made to classification performance, given the informal notion that it is easier to buy more memory than it is to attain higher classification accuracy.

This argument is made clearer when we evaluate the impact of a reduced feature space on the proto zero-day testing set. To this end, sub-Fig. 3B shows that the best performing model using 100% of features, in terms of accuracy, was the CNN v2 with a score of 0.842. When this model was trained and tested on the set of 30% of features, the accuracy declined to 0.751. Furthermore, sub-Fig. 3G shows that the FPR increased from 0.038 to 0.076, and similarly that the FNR declined from 0.255 to 0.389, evident in sub-Fig. 3H. An explanation for the more significant decline in classification performance when evaluated on the proto zero-day set compared to the normal testing set is that detection of botnet traffic from a range of known botnet families is a fairly simple task, enabling (relatively) high accuracies to be obtained with fewer features. On the other hand, echoing the discussion in Sect. 6.1, detection of botnet traffic from a range of unknown botnet families requires models to learn complex patterns shared by all botnets, which the datasets with reduced feature spaces are not rich enough to support.

## 6.2 Multiclass Classifiers for Botnet Classification

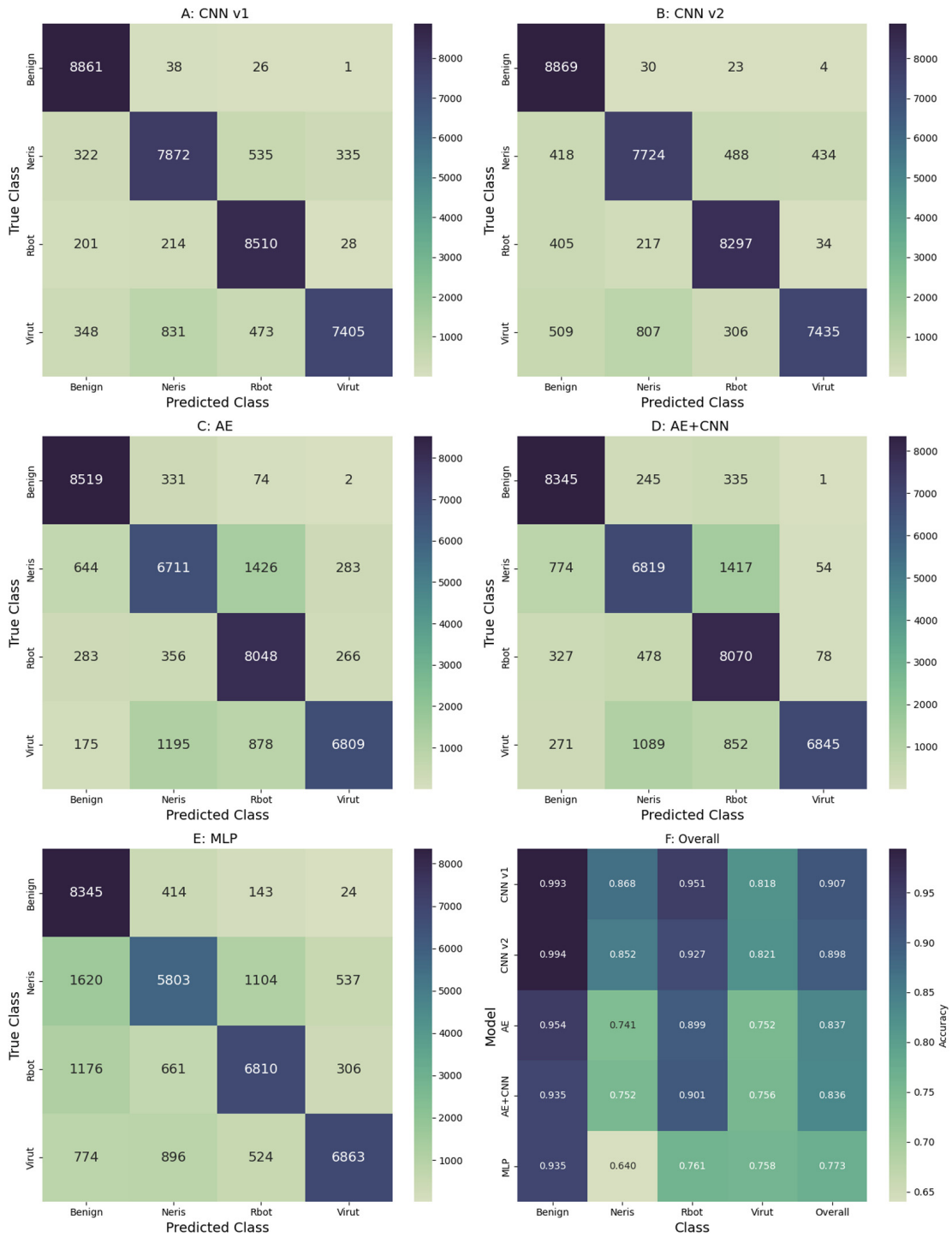
**Classification Accuracies of Each Model Overall, and for Each Botnet Family.** The classification performance of each of the five multiclass classifiers is shown in Fig. 4; each model has a corresponding confusion matrix which displays its predictions, and enables evaluation of how the model performs w.r.t. each botnet family. Sub-Figure 4F offers a holistic representation of all the models' performance on each respective family. Evident in Fig. 4F, the overall accuracies of the models when classifying traffic into respective families of Benign, Neris, Rbot, and Virut were expectedly lower than the binary classification task model accuracies, with a mean accuracy score across all classes of 0.850, compared to 0.986. The best performance was observed in the CNN v1, which achieved an average accuracy across all families of 0.907. We observed that there was again a marked improvement, in terms of overall classification accuracy, seen in the models which used convolutional layers, with the exception of the AE+CNN model. CNNs are known for their efficacy when learning hierarchical relationships between features, which is a useful way of reasoning about the necessary and sufficient conditions when making a classification [6]. In the context of this classification task, this ability may be an explanation for their performance, as

they are better able to combine the surface level and more complex features of network flows in order to learn more detailed patterns from the data.

As with the binary classification task, the shallow CNN (v1) outperforms the deep CNN (v2) by a small margin, with accuracies of 0.907 and 0.898, respectively. This is a percentage point increase of 0.009, which is fairly negligible. These findings allow us to make the claim that whichever patterns are learned by the CNNs are able to be learned with a shallow network, and that additional complexity (and associated dropout layers to combat overfitting) is unnecessary.

In the discussion from Sect. 6.1, concerning binary classification accuracy on the proto zero-day set, we suggested that the models that implemented convolutional layers (CNN v1, CNN v2, and AE+CNN) performed better due to the ability of CNNs to learn complex underlying behaviour patterns present in all botnet families. With respect to multiclass classification accuracy, CNN v1 and v2 are the best performers with overall accuracies of 0.907 and 0.898, respectively; while the AE+CNN has substantially worse classification performance, with an accuracy of 0.837. We have just proposed that a CNNs' ability to capture hierarchical relationships in data may be an explanation for their effectiveness for this problem; consequently, we suggest that an explanation for the relatively poor performance of the AE+CNN, which should benefit from this property of the convolutional layers, is that the encoding phase of the network reduces the complexity of the data to a point where the subsequent CNN is unable to learn the necessary hierarchical relationships, because they no longer 'exist' in the encoded representation. We further suggest that the reason for this is that the typical purpose of an autoencoder is to encode the input into a lower-dimensionality representation, from which an approximation of the original input can be reconstructed. This process may encourage the encoder to prioritise the more 'visible', surface-level patterns as they would be the best way to approximate the input. The consequence being that the encoder begins to act as a bottleneck - the features which aid more complex pattern learning are not present in the encoding, preventing a subsequent CNN from exploiting them. When we compare the results of the AE+CNN to the AE, similar accuracy rates are observed across the board, with overall accuracies of 0.837 and 0.836, respectively. This appears to reaffirm the notion of the encoder acting as a bottleneck for further classification.

Marín et al. [12] implemented a multiclass classifier of a deep CNN fed into an LSTM, trained on a dataset of Benign, Neris, Rbot, and Virut classes. Our best performing model, the CNN v1, showed an improvement on their CNN+LSTM, with an overall accuracy of 0.907 to their 0.765. Apart from the model architecture, a significant difference between our model and theirs is that our CNN v1 uses network flows as data, while they use bytes from the packet payloads, which are encrypted. We suggest that the difference in classification performance is caused by their model struggling to learn meaningful representations from the encrypted data. In support of this notion we refer to their binary classification experiments, where their payload-based classifier obtained an accuracy of 0.650, a result significantly lower than their flow-based classifier at 0.900, and our best performing binary classifier at 0.979 [12].



**Fig. 4.** Sub-Figures A to E show the respective confusion matrices for CNN v1, CNN v2, AE, AE+CNN, and MLP models, showing performance of each model w.r.t. individual classes. Sub-Figure F shows the overall accuracy, false-positive rate, and false negative rate of each model w.r.t. each class.

Notably, the CNN+LSTM struggled the most when classifying instances of Neris and Virut families. Figures 4B and F quite clearly show that this trend is observable across all models. Marín et al. suggest that an explanation of this is

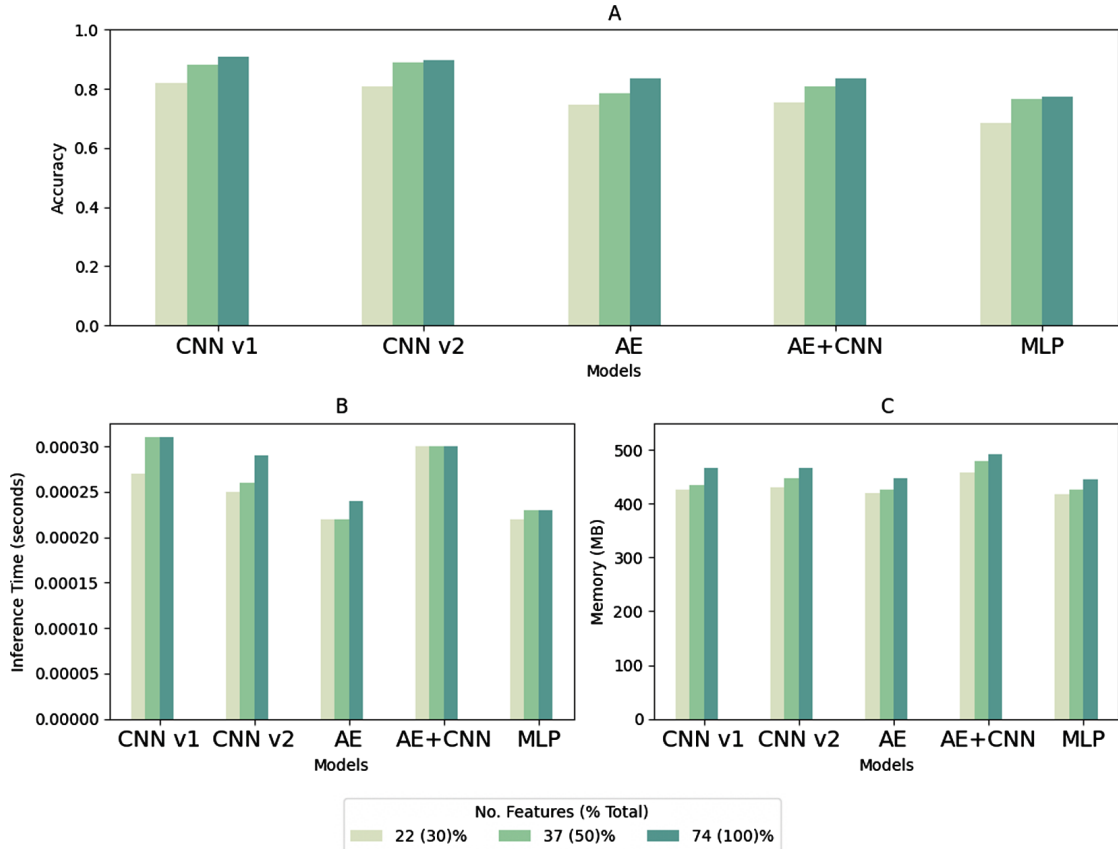


a result of the similarity between the Neris and Virut botnet families, causing models to misclassify one as the other. The results in Fig. 4 show that Virut samples are most frequently misclassified as Neris, supporting this notion. However, Neris samples are most frequently misclassified as Rbot, which may indicate that there is some other cause for the models’ confusion.

**Impact of Reduced Feature Space on Computational and Multiclass Classification Performance.** Figure 5A shows the accuracy of each model given the size of the feature space. As the number of features is reduced we observe an associated decline in classification accuracy. This aligns with the intuition that more features enable a model to use the relationships between features to learn more complex patterns, facilitating better classification performance. The two CNN models show a more substantial decline in accuracy when the features are reduced from 50% to 30%, when compared to 100% reduced to 50%. In some sense, this is an unintuitive result; the larger reduction in feature space is accompanied by a smaller reduction in accuracy. An explanation for this lies in the notion that the CNN’s success is a consequence of their ability to learn useful patterns from relations between features. When the feature space is reduced from 100% to 50%, there remains a sufficient number of features to enable the models to learn these patterns. In the reduction from 50% to 30%, while fewer features are removed, the resultant dataset is not detailed enough for CNN’s to learn useful information. However, an alternate explanation of why the decline in accuracy is more pronounced when going from 50% to 30% of feature space, as opposed to 100% to 50%, is that the process of reducing the feature space uses random sampling to select features. This may have resulted in important features being absent from the 30% dataset. To this end, the cause of the decline might be a result of quality, rather than quantity of features.

We observe that the three models that employ a CNN appeared to use more memory than the AE and MLP. We expect the MLP to use the least memory, as it is the least complex model. However, the CNN v1 has fewer trainable parameters than the AE, while using more memory. We suggest that the cause of this, and a general explanation for why CNNs seem to have the largest MMU, is that the CNN has to store filters and their respective activation maps in memory, which can become fairly expensive [6].

We find that in every model, there is an increase to MMU as the feature space increases. This was an expected result, reaffirming the position outlined in Sect. 6.1 that larger input feature spaces are associated with an increase to a model’s MMU. However, we maintain that this increase to MMU represents a relatively small improvement to computational performance, and is often coupled with a fairly substantial improvement in classification performance. For instance, when going from 30% to 100% of features the MLP’s MMU performance declines, with an increased utilisation of only  $\approx 28$  MB; however, there is an accompanied improvement to accuracy of 0.089.



**Fig. 5.** Sub-Figures A, B, and C show the respective accuracy, Mean Inference Time, and Mean Memory Usage when evaluated on datasets of varying feature size.

The CNN v1, v2 and AE+CNN models have the three slowest inference times when evaluated on 100% of the feature space. This result matches our expectations, as CNNs operate by executing a convolution operation for each filter across all output elements of the preceding layer [6]. This convolutional process is computationally intensive and is not a requirement in the Autoencoder (AE) and Multilayer Perceptron (MLP) architectures. While these trends seemingly continued as the feature space was reduced, we found no discernible relation between the MIT and the size of the feature space.

## 7 Conclusions and Future Work

In this paper we evaluated the classification and computational performance of DL models for botnet detection, and classification. We further explored the effects of reducing feature spaces on performance. For our first research objective, we found that all models achieved accuracy  $\geq 0.979$ , FPRs  $\leq 0.033$ , and FNRs  $\leq 0.026$ , suggesting that the classification problem was fairly simple. Classification performance on the proto zero-day set saw models which implemented convolutional layers have substantially better classification performance. Additionally, there was considerably more variation between models' performance, indicating that architecture choice plays a more significant role when detecting

unknown botnet families. We suggest this was a result of the ability of convolutional layers to learn the more complex behaviour, general to botnets, as opposed to surface-level features. Learning the behaviour patterns of botnets in general enables the models to perform well on testing sets made up of entirely unknown botnet families. We further observed a clear trend that larger feature spaces were associated with a larger memory utilisation (MMU), affirming our expectation that feature selection might improve computational performance. However, in reducing feature space we also observe a substantial decline in classification performance. We found no evidence of a trend between feature space size and MIT; however, we acknowledge that there were serious limitations to the accuracy of measuring MIT. The consequence of this is that an 'optimal' size for a feature space should be determined on a per case basis, it is dependent on relevant constraints.

While we found little evidence of a reduced feature space improving the inference time of models, we did observe that CNNs in general seemed to have longer inference times, which we attributed to the expense of the convolution operation. Furthermore, we found that CNNs appeared to have a greater associated MMU, which we suggest is a result of the filters and respective activation maps, and not solely their complexity. With respect to MMU across varying feature space, we observed in all models that a reduction in feature space was clearly associated with a lower MMU. However, as was argued with respect to the binary classification problem, the fairly minor improvements to the computational cost of models are accompanied by an arguably more significant degradation to classification performance.

With respect to the second research objective, we found a fairly large differential between the classification performance of models which used convolutional layers and those that did not. We suggest that the efficacy of CNNs at learning hierarchical relations between features is an explanation for this. As with the binary classifiers, we observed a trend where larger feature spaces were associated with slightly greater MMU. However, the larger feature spaces resulted in substantially better classification performance.

The most significant observation from this work is in the classification performance of binary classifiers on the proto zero-day set. In this area, we observed accuracies high enough to act as a proof of concept: that models are able to generalise what they learn from training data to entirely unseen botnets. However, these results leave considerable room for improvement.

Future work might take one of two directions; while we explored the theoretical capabilities of DL models to detect and classify encrypted botnet traffic, an interesting extension would be a practical implementation of some of this work. We suggest that a good starting point would be to implement a binary classifier for detection - as detection, when contrasted with family classification, seems to have more immediate utility. On the theoretical side, we suggest ablation studies to investigate both feature and hyperparameter importance. The notion being that model size could be reduced if we knew more about the importance of certain features, or parts of the model. This may further yield insight into

behaviours common between botnets, and consequently enable tuning models specifically for handling new botnet families.

## A Appendix

For comprehensive information, see [Project Website](#) or [Github](#).

### A.1 Feature Sets

**Table 2.** All the features present in flow extraction from CICFlowMeter. Bold and/or underlined indicate inclusion in 30% and 50% feature-spaces, respectively. All features present in 100% feature space.

No.	Feature Name	No.	Feature Name	No.	Feature Name
1	Flow ID	30	<b><u>Fwd IAT Max</u></b>	59	<b>Average Pkt Size</b>
2	Src IP	31	Fwd IAT Min	60	Fwd Segment Size Avg
3	Src Port	32	<u>Bwd IAT Total</u>	61	Bwd Segment Size Avg
4	Dst IP	33	<u>Bwd IAT Mean</u>	62	<b>Fwd Bytes/Bulk Avg</b>
5	Dst Port	34	<b>Bwd IAT Std</b>	63	<b>Fwd Pkt/Bulk Avg</b>
6	Protocol	35	<u>Bwd IAT Max</u>	64	Fwd Bulk Rate Avg
7	Timestamp	36	Bwd IAT Min	65	Bwd Bytes/Bulk Avg
8	<u>Flow Duration</u>	37	<u>Fwd PSH Flags</u>	66	<u>Bwd Pkt/Bulk Avg</u>
9	<u>Total Fwd Pkt</u>	38	Bwd PSH Flags	67	<u>Bwd Bulk Rate Avg</u>
10	<u>Total Bwd Pkts</u>	39	Fwd URG Flags	68	<u>Subflow Fwd Pkts</u>
11	<b><u>Total Length of Fwd Pkt</u></b>	40	Bwd URG Flags	69	<b>Subflow Fwd Bytes</b>
12	<b><u>Total Length of Bwd Pkt</u></b>	41	<b>Fwd Header Length</b>	70	Subflow Bwd Pkts
13	<b><u>Fwd Pkt Length Max</u></b>	42	<u>Bwd Header Length</u>	71	<b>Subflow Bwd Bytes</b>
14	Fwd Pkt Length Min	43	<b>Fwd Pkts/s</b>	72	<u>FWD Init Win Bytes</u>
15	<u>Fwd Pkt Length Mean</u>	44	<u>Bwd Pkts/s</u>	73	<u>Bwd Init Win Bytes</u>
16	<u>Fwd Pkt Length Std</u>	45	<b>Pkt Length Min</b>	74	<u>Fwd Act Data Pkts</u>
17	Bwd Pkt Length Max	46	<u>Pkt Length Max</u>	75	<b>Fwd Seg Size Min</b>
18	Bwd Pkt Length Min	47	<u>Pkt Length Mean</u>	76	<u>Active Mean</u>
19	Bwd Pkt Length Mean	48	Pkt Length Std	77	<b>Active Std</b>
20	<b><u>Bwd Pkt Length Std</u></b>	49	<u>Pkt Length Variance</u>	78	Active Max
21	Flow Bytes/s	50	FIN Flag Count	79	<b>Active Min</b>
22	Flow Pkts/s	51	<u>SYN Flag Count</u>	80	Idle Mean
23	Flow IAT Mean	52	RST Flag Count	81	<b>Idle Std</b>
24	<u>Flow IAT Std</u>	53	PSH Flag Count	82	<u>Idle Max</u>
25	<u>Flow IAT Max</u>	54	<u>ACK Flag Count</u>	83	<b>Idle Min</b>
26	Flow IAT Min	55	<b><u>URG Flag Count</u></b>	84	<b>Label</b>
27	Fwd IAT Total	56	CWR Flag Count		
28	<u>Fwd IAT Mean</u>	57	<b>ECE Flag Count</b>		
29	<u>Fwd IAT Std</u>	58	Down/Up Ratio		

## A.2 Classification Results

(See Tables 3, 4, 5 and 6)

**Table 3.** Results of Binary Classifiers on Default Test Set

Model	Accuracy			Precision			Recall			FPR			FNR		
	100	50	30	100	50	30	100	50	30	100	50	30	100	50	30
CNN v1	.990	.972	.966	.990	.972	.966	.990	.972	.966	.016	.030	.032	.004	.026	.036
CNN v2	.993	.983	.968	.993	.983	.968	.993	.983	.968	.008	.021	.027	.006	.012	.037
AE	.979	.974	.966	.979	.974	.966	.979	.974	.966	.030	.026	.035	.012	.027	.033
AE CNN	.983	.974	.963	.983	.974	.963	.983	.974	.963	.008	.030	.034	.026	0.022	.040
MLP	.986	.971	.959	.986	.971	.959	.986	.971	.959	.008	.033	.055	.021	0.026	.027

**Table 4.** Results of Binary Classifiers on Proto Zero-Day Test Set

Model	Accuracy			Precision			Recall			FPR			FNR		
	100	50	30	100	50	30	100	50	30	100	50	30	100	50	30
CNN v1	.793	.653	.480	.904	.777	.671	.700	.736	.116	.092	.448	.071	.300	.264	.884
CNN v2	.842	.714	.751	.960	.777	.909	.745	.677	.611	.038	.240	.076	.255	.323	.389
AE	.732	.552	.716	.783	.581	.824	.713	.683	.619	.245	.609	.163	.287	.317	.381
AE CNN	.810	.705	.687	.957	.727	.785	.688	.746	.597	.038	.346	.202	.312	.254	.403
MLP	.740	.594	.698	.782	.610	.744	.734	.736	.691	.253	.580	.294	.266	.264	.309

**Table 5.** Computational Performance of Binary Classifiers

Model	Memory (MB)			Inference Time (S)		
	100	50	30	100	50	30
CNN v1	465.89	422.91	422.39	0.00032	0.00024	0.00027
CNN v2	526.16	443.78	435.83	0.00031	0.00026	0.00026
AE	457.02	428.02	420.16	0.00022	0.00022	0.00023
AE CNN	490.94	485.50	475.42	0.00026	0.00023	0.00033
MLP	446.05	422.70	416.67	0.00020	0.00020	0.00020

**Table 6.** Classification and Computational Performance of Multiclass Classifiers

Model	Accuracy (%)			Memory (MB)			Inference Time (S)		
	100	50	30	100	50	30	100	50	30
CNN v1	.907	.881	.819	465.96	434.63	425.77	.00030	.00031	.00027
CNN v2	.898	.889	.809	465.85	446.88	431.28	.00029	.00026	.00025
AE	.836	.786	.745	447.98	427.21	419.70	.00024	.00022	.00022
AE CNN	.836	.807	.755	492.95	478.88	458.87	.00028	.00030	.00030
MLP	.773	.764	.684	445.69	425.44	418.13	.00023	.00023	.00022

## References

1. Abu Rajab, M., Zarfoss, J., Monrose, F., Terzis, A.: A multifaceted approach to understanding the botnet phenomenon. In: Proceedings of the 6th ACM SIGCOMM Conference on Internet Measurement, IMC 2006, pp. 41–52. Association for Computing Machinery, New York (2006). <https://doi.org/10.1145/1177080.1177086>
2. Aceto, G., Ciunzo, D., Montieri, A., Pescapé, A.: Mobile encrypted traffic classification using deep learning. In: 2018 Network Traffic Measurement and Analysis Conference (TMA), pp. 1–8. IEEE (2018)
3. Bertino, E., Islam, N.: Botnets and internet of things security. *Computer* **50**(2), 76–79 (2017)
4. Cheng, R.: D 2 pi : identifying malware through deep packet inspection with deep learning (2017). <https://api.semanticscholar.org/CorpusID:53062187>
5. García, S., Grill, M., Stiborek, J., Zunino, A.: An empirical comparison of botnet detection methods. *Comput. Secur.* **45**, 100–123 (2014). <https://doi.org/10.1016/j.cose.2014.05.011>, <https://www.sciencedirect.com/science/article/pii/S0167404814000923>
6. Goodfellow, I., Bengio, Y., Courville, A.: Deep Learning. MIT Press (2016). <http://www.deeplearningbook.org>
7. Haddadi, F., Le Cong, D., Porter, L., Zincir-Heywood, A.N.: On the effectiveness of different botnet detection approaches. In: Lopez, J., Wu, Y. (eds.) ISPEC 2015. LNCS, vol. 9065, pp. 121–135. Springer, Cham (2015). [https://doi.org/10.1007/978-3-319-17533-1\\_9](https://doi.org/10.1007/978-3-319-17533-1_9)
8. Lashkari, A.H., Gil, G.D., Mamun, M.S.I., Ghorbani, A.A.: Characterization of tor traffic using time based features. In: Proceedings of the 3rd International Conference on Information Systems Security and Privacy - Volume 1: ICISSP, pp. 253–262. INSTICC, SciTePress (2017). <https://doi.org/10.5220/0006105602530262>
9. Li, L., Jamieson, K., DeSalvo, G., Rostamizadeh, A., Talwalkar, A.: Hyperband: a novel bandit-based approach to hyperparameter optimization. *J. Mach. Learn. Res.* **18**(1), 6765–6816 (2017)
10. Lim, H.K., Kim, J.B., Kim, K., Hong, Y.G., Han, Y.H.: Payload-based traffic classification using multi-layer LSTM in software defined networks. *Appl. Sci.* **9**(12), 2550 (2019)
11. Lotfollahi, M., Jafari Siavoshani, M., Shirali Hossein Zade, R., Saberian, M.: Deep packet: a novel approach for encrypted traffic classification using deep learning. *Soft Comput.* **24**(3), 1999–2012 (2020)
12. Marín, G., Caasas, P., Capdehourat, G.: DeepMAL - deep learning models for malware traffic detection and classification. In: Data Science – Analytics and Applications, pp. 105–112. Springer, Wiesbaden (2021). [https://doi.org/10.1007/978-3-658-32182-6\\_16](https://doi.org/10.1007/978-3-658-32182-6_16)
13. O’Malley, T., et al.: Kerastuner (2019). <https://github.com/keras-team/keras-tuner>
14. O’Shea, K., Nash, R.: An introduction to convolutional neural networks. arXiv preprint [arXiv:1511.08458](https://arxiv.org/abs/1511.08458) (2015)
15. Pachhala, N., Jothilakshmi, S., Battula, B.P.: A comprehensive survey on identification of malware types and malware classification using machine learning techniques. In: 2021 2nd International Conference on Smart Electronics and Communication (ICOSEC), pp. 1207–1214 (2021). <https://doi.org/10.1109/ICOSEC51865.2021.9591763>



16. Papadogiannaki, E., Tsirantonakis, G., Ioannidis, S.: Network intrusion detection in encrypted traffic. In: 2022 IEEE Conference on Dependable and Secure Computing (DSC), pp. 1–8 (2022). <https://doi.org/10.1109/DSC54232.2022.9888942>
17. Acarman, T.: Botnet detection based on network flow summary and deep learning. *Int. J. Netw. Manage.* **28**(6), e2039 (2018). <https://doi.org/10.1002/nem.2039>, <https://onlinelibrary.wiley.com/doi/abs/10.1002/nem.2039>
18. Piskozub, M., Gaspari, F.D., Barr-Smith, F., Mancini, L., Martinovic, I.: Mal-Phase: fine-grained malware detection using network flow data. In: Proceedings of the 2021 ACM Asia Conference on Computer and Communications Security. ACM (2021). <https://doi.org/10.1145/3433210.3453101>
19. van Roosmalen, J., Vranken, H., van Eekelen, M.: Applying deep learning on packet flows for botnet detection. In: Proceedings of the 33rd Annual ACM Symposium on Applied Computing, pp. 1629–1636 (2018)
20. Sarker, I.H.: Cyberlearning: Effectiveness analysis of machine learning security modeling to detect cyber-anomalies and multi-attacks. *Internet Things* **14**, 100393 (2021)
21. Stratosphere: Stratosphere laboratory datasets (2015). <https://www.stratosphereips.org/datasets-overview>. Accessed 13 Mar 2020
22. Torres, P., Catania, C., Garcia, S., Garino, C.G.: An analysis of recurrent neural networks for botnet detection behavior. In: 2016 IEEE Biennial Congress of Argentina (ARGENCON), pp. 1–6. IEEE (2016)
23. Van Rossum, G., Drake, F.L.: Python 3 Reference Manual. CreateSpace, Scotts Valley (2009)
24. Villa, A., Varki, E.: Characterization of a campus internet workload. In: Proceedings of CATA, pp. 140–148 (2012)
25. Wang, W., et al.: HAST-IDS: learning hierarchical spatial-temporal features using deep neural networks to improve intrusion detection. *IEEE Access* **6**, 1792–1806 (2017)
26. Wang, Z., Fok, K.W., Thing, V.L.: Machine learning for encrypted malicious traffic detection: approaches, datasets and comparative study. *Comput. Secur.* **113**, 102542 (2022). <https://doi.org/10.1016/j.cose.2021.102542>
27. Weisz, S., Chavula, J.: Community network traffic classification using two-dimensional convolutional neural networks. In: Sheikh, Y.H., Rai, I.A., Bakar, A.D. (eds.) AFRICOMM 2021. LNICST, pp. 128–148. Springer, Cham (2022). [https://doi.org/10.1007/978-3-031-06374-9\\_9](https://doi.org/10.1007/978-3-031-06374-9_9)
28. Yeo, M., et al.: Flow-based malware detection using convolutional neural network. In: 2018 International Conference on Information Networking (ICOIN), pp. 910–913 (2018). <https://doi.org/10.1109/ICOIN.2018.8343255>
29. Zeng, Y., Gu, H., Wei, W., Guo, Y.: *deep – full – range*: a deep learning based network encrypted traffic classification and intrusion detection framework. *IEEE Access* **7**, 45182–45190 (2019). <https://doi.org/10.1109/ACCESS.2019.2908225>
30. Zhou, H., Hu, Y., Yang, X., Pan, H., Guo, W., Zou, C.C.: A worm detection system based on deep learning. *IEEE Access* **8**, 205444–205454 (2020)