

A Machine Learning-based Tool for Passive OS Fingerprinting with TCP Flavor as a Novel Feature

Destá Haileselassie Hagos, *Student Member, IEEE*, Anis Yazidi, *Senior Member, IEEE*, Øivind Kure, *Senior Member, IEEE*, and Paal E. Engelstad, *Senior Member, IEEE*

Abstract—With the emergence of Internet of Things (IoT), securing, and managing large, complex enterprise network infrastructure requires capturing and analyzing network traffic traces in real-time. An accurate passive Operating System (OS) fingerprinting plays a critical role in effective network management and cybersecurity protection. Passive fingerprinting doesn't send probes that introduce extra load to the network and hence it has a clear advantage over active fingerprinting since it also reduces the risk of triggering false alarms. This paper proposes and evaluates an advanced classification approach to passive OS fingerprinting by leveraging *state-of-the-art* classical machine learning and deep learning techniques. Our controlled experiments on benchmark data, emulated and realistic traffic is performed using two approaches. Through an Oracle-based machine learning approach, we found that the underlying TCP variant is an important feature for predicting the remote OS. Based on this observation, we develop a sophisticated tool for OS fingerprinting that first predicts the TCP flavor using passive traffic traces and then uses this prediction as an input feature for another machine learning algorithm for predicting the remote OS from passive measurements. This paper takes the passive fingerprinting problem one step further by introducing the underlying predicted TCP variant as a distinguishing feature. In terms of accuracy, we empirically demonstrate that accurately predicting the TCP variant has the potential to boost the evaluation performance from 84% to 94% on average across all our validation scenarios and across different types of traffic sources. We also demonstrate a practical example of this potential, by increasing the performance to 91.2% and 95.3% on average using a tool for loss-based and delay-based TCP variants prediction in an emulated setting. To the best of our knowledge, this is the first study that explores the potential for using the knowledge of the TCP variant to significantly boost the accuracy of passive OS fingerprinting.

Keywords—*Operating System, Fingerprinting, Machine Learning, Deep Learning, IoT, Passive Traffic Measurements*

I. INTRODUCTION AND MOTIVATION

AS modern network infrastructures grow in size, collecting detailed relevant knowledge about the dynamic characteristics and complexity of large heterogeneous networks

is crucial for many purposes e.g., network vulnerability assessment and monitoring, spam detection, etc. The interconnection and heterogeneity of IoT-enabled devices connected to the Internet also raises potential security issues and it has gained a lot of research attention from the industry to academia [3, 25, 45, 54]. Developing advanced network security and monitoring techniques are important for both the research and security practitioners. There has been a significant research work in the context of network management and cybersecurity on developing network security tools to fingerprint remote and local Operating Systems (OSes) [34, 35, 36, 55, 56]. OS fingerprinting is the process of inferring the information about the underlying OS of a machine operating with TCP/IP packets by a remote device connected on the Internet without having physical access to the device [27].

There are many different existing custom tools for fingerprinting of the most commonly used OSes based on the characteristics of its underlying TCP/IP network stack [27] and this, to a large extent, is due to variability in how the TCP/IP stack is traditionally implemented across different variations of OSes [33]. One common approach, for example, is by collecting the TCP/IP stack basic parameters [31], e.g., IP initial Time To Live (TTL) default values [8], HTTP packets using the User-agent field [30], Internet Control Message Protocol (ICMP) requests [39], known open port patterns, the size of the TCP receiver window [24], TCP Maximum Segment Size (MSS) [41], IP Don't Fragment (DF) flag [40], a set of other specific TCP options to mention a few. However, in our work, we want to take this one step further by combining these basic features and other settings with the underlying TCP variant as a feature in our model due to the fact that different OSes have slightly different implementations of TCP. Some implementations of common TCP variants quickly overshoot the size of the Congestion Window (cwnd) because of differences in the variant implementations. Hence, we believe that knowing the implementation of the underlying OS may help us understand better their exact behavior. It can also help us explore how to classify an OS when different classes of OSes are implementing the same TCP variant.

Fingerprinting Techniques: We can determine what OS a remote computer on the Internet is running by either passively listening to traffic captured from a network or by actively sending it packets. The most widely used complementary remote OS fingerprinting proven approaches that employ a variety of TCP/IP stack scanning are broadly categorized into classes of *active* and *passive* techniques.

D. Hagos, Ø.Kure, and P. Engelstad are with the Autonomous Systems and Sensor Technologies Research Group, Department of Technology Systems, University of Oslo, Oslo, 0316, Norway (e-mail: destahh@ifi.uio.no; oivind.kure@its.uio.no; paal.engelstad@its.uio.no).

D. Hagos, A.Yazidi, and P. Engelstad are with the Autonomous System and Network (ASN) Research Group, Department of Computer Science, Oslo Metropolitan University, Oslo, 0130, Norway (e-mail: desta.hagos@oslomet.no; anis.yazidi@oslomet.no; paalen@oslomet.no).

Submitted For Publication in IEEE Internet of Things Journal

- **Active Fingerprinting:** This technique is based on actively transmitting one or more specially crafted network packets with different packet settings or flags to a remote network device in order to analyze the corresponding potentially identifying replies [34, 55]. This method determines knowledge of the underlying OS according to the received responses from the target device by examining the network behavior of known TCP/IP stack [46]. However, since this approach injects additional traffic to the network by generating active probes, it may itself trigger alarms and get blocked by firewall rules and Network address translators (NATs) [14].
- **Passive Fingerprinting:** This approach, on the other hand, inspects and analyzes packets traveling between end hosts without injecting any traffic into the network [35, 36, 56]. This technique with little resource simply analyzes a pattern of the OS-specific information that has already been sent in the network traffic and compares for a match with a predefined database that contains a list of known signatures of different OSes. Passive fingerprinting doesn't send probes and hence it has a clear advantage over active fingerprinting since it reduces the risk of triggering alarms [14].

OS fingerprinting can also be performed using classical techniques known as “*banner grabbing*”. It is an approach used to gain detailed information about a remote computer system on a network and the associated services running on its opened ports [43]. Using techniques like this, some remote computers announce their underlying OS freely and running application services with their versions in use to anyone connecting to them as part of welcome banners or header information. Some of the widely used services that serve *banner grabbing* are: *Telnet*, *FTP*, *NetCat*, *SMTP*, etc. However, it is useful to remember that some of these basic services are effective against less secure networks.

Potential benefits and applications: Network scanning and accurate remote OS fingerprinting are the crucial steps for penetration testing in terms of security and privacy protection. Note that attackers can also embrace passive fingerprinting techniques to search for potential victims in a network. For example, by identifying the OS running on a remote computer and the list of services it runs, an attacker can target the device to eavesdrop on the communication between the endpoints without having physical access to the device. However, we argue that our work presented here is motivated by a number of practical applications that can be positively used by network and system administrators. Passively fingerprinting an OS by analyzing the packets it generates and transmits over a network is extremely important in the areas of network management and computer security for several reasons. For example, it is useful to explore a network for potential exploitations of security vulnerabilities which can be exploited by attackers, auditing, identify critical attacks, reveal new information about a network user etc. Network administrators can, therefore, use this OS related information to maintain the security policy and reliability of their network by configuring a network-based Intrusion Detection Systems (IDS) [32]. Vulnerabilities and

security threats in a network may result from rogue or unauthorized devices [49], unsecured internal nodes within the network, and from external nodes [7]. Hence, passively fingerprinting an OS has a potential benefit in addressing these critical problems. This, from an academic point of view, is interesting and something that needs to be addressed from a network security research point of view.

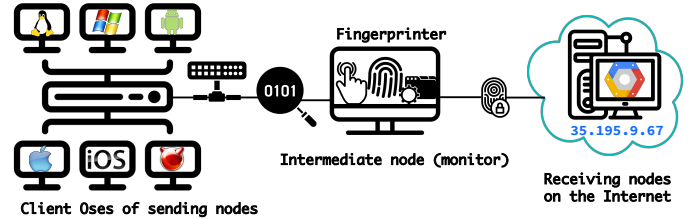


Fig. 1: Network architecture for passive OS fingerprinting by an intermediate node.

Limitations of previous works: Traditionally, most of the existing general OS fingerprinting techniques resort to manually generated signature matching from a database of heuristics which contains features of widely used OSes. This means, after comparing the generated signatures, the first set of responses match with the highest confidence against a database of fingerprints would be used to select the specific probable OS. However, manually updating a large number of signature and managing databases of new OSes adds a considerable amount of time and hence we may suffer from the consequences of the lack of recent signature updates of the known OSes. For example as reported in [30], the last updates of the fingerprint databases of *Ettercap* [36] and *p0f* [56] date to 2011 and 2014 respectively. Consequently, new OSes families like Android 4.4 and higher versions of Android, Windows 10 distributions, etc. will not be recognized by these tools since they are not included in their fingerprint databases. Hence, we argue that it is important to consider making use of a fingerprint database that contains variations of most currently used OSes and automating these tasks by employing learning algorithms capable of extracting all possible OS-specific features for discovering the underlying OSes. To explore this idea of applying machine learning algorithms, we present a unified and robust classification approach to an advanced passive OS fingerprinting that leverages both machine learning and deep learning methods. Our fingerprinting technique is completely passive meaning that we only need to be able to observe network traffic from a target machine at any observation point on the network without injecting any traffic into the network. Note that the TCP/IP header fields would not be impacted by SSL/TLS encryption of the TCP payload. Hence, since we utilize features that are readable even with encryption, our approach is independent of whether the flow is encrypted or not. Figure 1 shows the architecture for implementing our fingerprinting methodology.

Why use machine learning approaches to perform OS fingerprinting? There are several limitations imposed by classical fingerprinting techniques. Passive OS fingerprinting

generally relies on recognizing the default values for various TCP/IP stack parameters. Traditional approaches to OS fingerprinting resort to TCP header parameters such as the initial SYN packet size, TTL in the IP header, and the size of the TCP receiver window of the first packet in a TCP session which can be easily manipulated by experienced users for Quality of Service (QoS) or by an adversary with malicious intent. Therefore, the traditional approaches can be easily misled by changes in TCP/IP input features limiting the accuracy of OS classification. If a user changes these parameters, the task of OS fingerprinting becomes much more challenging. Most of the existing works on OS fingerprinting provide little capability to address this challenge.

Motivated by this problem, we proposed a novel approach by leveraging both machine learning and deep learning-based techniques that consider the set of parameters as a whole, rather than individually so that our model caters for variations in TCP parameters. If a user changes the initial receive window size, for instance, we may still be able to recognize the OS from other parameters that have not been changed (TCP congestion control algorithm, initial cwnd size, etc.). Note that this depends entirely on the changes made by the user to the default TCP or OS stack parameters that are commonly used for signature-based fingerprinting. In this paper, we investigate the potential of knowing the underlying TCP variant and how much it might improve the OS fingerprinting accuracy. The advantage of introducing the TCP variant as input feature is the fact that it is a characterizing feature of an OS that is difficult to manipulate. Modifying the TCP behavior of an operating system is not an easy task and needs changes to the kernel of the underlying operating system. Thus, TCP behavior is believed to be a robust input feature that is less prone to changes in the configuration by the user. The other reason why we create a model by employing learning techniques is to understand the complex patterns of the varying values in the TCP header and extract useful input features. Because machine learning offers new possibilities as it can extract patterns and general rules for classification. Machine learning can also be more robust to small variations in the input parameters. In addition to this, with the use of learning techniques, we argue that avoiding using manually updated static signature databases has two potential benefits. Firstly there is no tedious task of creating these unique fingerprints, all you need is a set of values or features. The second benefit comes from a known flaw in many of the existing fingerprinting tools, where a “first-match” policy is applied, meaning that if two fingerprints are equal the tool would always predict the first OS with that exact fingerprint. However, learning techniques, on the other hand, make calculated guesses of which of the classes with the same fingerprint that will be predicted.

Contributions: We summarize our main contributions below.

- We propose and evaluate a robust approach to OS fingerprinting from passive measurements by leveraging machine learning and deep learning techniques.
- We investigate the use of TCP congestion control variant as a distinguishing feature in passive OS fingerprinting.

- We explore variability in implementations of TCP variant by different OSes and its effect on classifying remote OS.
- We study the applicability of Recurrent Neural Networks (RNN)-based models for robust and advanced passive OS fingerprinting by combining the basic TCP/IP features and the predicted TCP variant as input vectors.
- We show that the TCP flavor has a great potential for boosting passive OS fingerprinting accuracy.
- We build a universal tool that can be applied to first estimate the TCP cwnd from passive measurements, second predict the underlying TCP flavor, and finally uses the predicted TCP variant as an input feature to fingerprint the remote computer’s OS.

Roadmap: The rest of the paper is organized as follows. Section II discusses related work, and Section III presents the experimental datasets. Section IV presents the machine learning of the OS fingerprinter. The role of TCP variant in passive OS fingerprinting and its feasibility across all use cases is presented in Section V. The experimental results with the loss-based and delay-based predicted TCP variants and the transfer learning are presented in Section VI. Finally, Section VII concludes our paper and suggests directions for future research work.

II. RELATED WORK

Remote OSes fingerprinting has a long history in the computer security community [2, 30, 31, 34]. In this section, we briefly summarize the relevant related works as follows.

TCP/IP header fingerprinting and any information related to application protocols are used to identify the underlying OS running on a remote host either actively or passively [33]. As we explained in Section I, there are multiple existing tools for both the predominant active and passive OS fingerprinting approaches, where *Nmap* [34] is one of the most prominent open-source active fingerprinting tools. The work presented in [47], *SYNSCAN*, works in a similar fashion to *Nmap*, but it performs the fingerprinting task by actively sending a small number of crafted network packets to a single TCP port. *Xprobe2* [55] is another popular remote active OS fingerprinting tool, which relies primarily on different types of ICMP packets. By actively sending a small number of User Datagram Protocol (UDP) and ICMP request packets to the remote target host, *Xprobe2* triggers ICMP datagram responses. *Xprobe2* uses a fuzzy logic matching algorithm based on a statistical computation of scores for each test performed. Since *Xprobe2* utilizes a simple fuzzy signature matching against the signature database, its detection accuracy is prone to small changes to the default TCP/IP stack parameters that might harden the detection of the underlying OS. However, *Xprobe2* is more robust to small fingerprint variations as compared to *Nmap*. This is due to the fact that *Nmap* uses static rule matching while *Xprobe2* enjoys more flexibility thanks to its fuzzy matching logic.

As explained above the other OS fingerprinting tools, *Ettercap* [36] and *p0f* [56], have not been updated since 2011 and 2014 respectively to include variations of most widely used modern OSes. For passive OS fingerprinting to be effective, we

believe that the limitations of these fingerprinting tools need to be addressed. The work in [31] also demonstrates that the OS fingerprinting accuracy of the *Ettercap* and *p0f* signature databases is low and techniques to improve performance was proposed. It presents rule-based machine learning classifiers capable of identifying 75 classes of OSes from TCP/IP packet headers found in the *Ettercap* database. Lippmann et al. proposed a classifier technique using k-nearest neighbors (KNN) that returns an approximate first match for an OS from a fingerprint database. This counters the problem of classifying remote and local hosts as unknown if no exact match is found in the database [31]. However, their evaluation yielded poor experimental results, rejecting as much as 84% of the test packets, while 44% of the accepted patterns were wrongly classified [31]. The problems contributing to poor OS classification performance were believed to be caused by two main reasons. The first reason is substitution errors due to multiple OSes sharing exactly the same fingerprint feature values. The second reason for this poor OS classification performance is the high rejection rate caused by numerous unique feature values derived from the same OS. This error can only be reduced by combining OS classes. After combining all the OS classes, the error percentage was reduced to 9.8% with no rejected packets. It is worth mentioning that fingerprinting techniques have been also extended to remote device level fingerprinting [12], thus going beyond remote OS detection using TCP/IP network stacks.

A recent study that is most closely related to our work, and which has also given a comprehensive survey on passive fingerprinting methods, can be found in [30]. The authors have employed OS fingerprinting methods in the environment of wireless networks. Besides using the three basic TCP/IP stacks (i.e., TTL, window size, and initial SYN packet size), the authors suggested also using the user-agent information in HTTP request headers and communication with OS-specific domains can be usable in large dynamic networks [30]. As shown in Table II, the average accuracy of OS classification using the TCP/IP parameters reported in [30] is 80.88%. Zhang et al.’s paper on OS detection [57] utilizes only one machine learning technique namely Support Vector Machine (SVM). However, the testing error rate of identifying some of the OSes e.g., Mac, Cisco, FreeBSD, and OpenBSD is 25.80%, 24.22%, 17.71%, and 15.85% respectively [57]. Aksoy et al. [2] have employed genetic algorithms for identifying packet features suitable for OS classification based on the analysis of the network TCP/IP packets using machine learning algorithms. However, most of these previous works use the basic actual TCP/IP features for evaluating passive OS fingerprinting. Besides, we believe that these tools have the inability to extract all possible OS-specific features for passively fingerprinting the underlying OSes. Examples of those features include OS-specific Domain Name System (DNS) queries [6], Dynamic Host Configuration Protocol (DHCP) options [28], etc.

In contrast, what separates our contribution in this paper from the other previous related works is that our tool supports a wider range of TCP/IP network stack features. As shown in Figure 2, the main goal of our work presented

here is to combine these basic TCP/IP features that are the basis of OS fingerprinting with the underlying TCP variant by leveraging both machine learning and deep learning techniques. This contribution remains largely unexplored and is not used by existing OS fingerprinting techniques. Detecting the implementation of a TCP variant passively is a challenging task and this, we believe, is the reason why no previous works use it to passively fingerprint remote and local OSes. However, in our case, we already have a general solution for this difficulty presented in our previous works [18, 19, 20]. The reason why we focus on the implementations of the underlying TCP variant as a feature in our OS classifier model is due to the fact that different OSes are doing slightly different implementations of TCP. Hence, we believe that passively observing the network-level characteristics found in TCP packets can give us more information about the remote computer’s underlying OS. We further believe that this will also help us to explore in detail the long-term characteristics of TCP traffic. To the best of our knowledge, this is the first study that sheds light on the potential of the underlying TCP variant feature in boosting significantly the accuracy of passive OS fingerprinting using machine learning and deep learning techniques. [The classification performance of our machine learning and deep learning approaches in comparison with other state-of-the-art tools for passive OS fingerprinting, such as p0f \[56\], is presented in Table XVI.](#)

III. EXPERIMENTAL DATASETS

Our machine learning models for OS classification is developed and tested on three datasets, presented below.

A. Benchmark Data

First, we utilize a large benchmark dataset that has been used for OS fingerprinting in a previous related work [30]. This dataset is closely aligned with our task, and it was collected from a university wireless network. The benchmark dataset was used in the previous work for OS fingerprinting based on the HTTP header, while the ambition of our paper is to do generic fingerprinting based only on the TCP packet fields. Since we aim at fingerprinting that is not application-specific, the TCP information in the dataset is useful for our purpose, while the HTTP User-agent information in our experiments is used only to establish ground truth about the OS that was used.

TABLE I: Statistics and distribution of the OSes and their market shares within the OS-family.

Android	Windows	Mac OS	Linux	iOS	Unix	Other
8.0	10	Mojave	Ubuntu 16.04	12.1	Solaris 11.4	Unknown
8.1	7	High Sierra	Ubuntu 18.04	11.4	FreeBSD 11.2	
6.0	8.1	Sierra	Ubuntu 18.10	12.0		
7.0	8.0	El Capitan	Fedora 29	10.3		
7.1	XP	Yosemite	Debian 9	9.3		
5.1	Vista	Mavericks	CentOS 7.6	11.2		
			openSUSE 42.3			
36.5%	35.99%	6.37%	0.79%	13.99%	1.58%	

The benchmark dataset contains 79087345 flows, activity of 21746 unique users, 253374 WiFi sessions, 25642 unique

MAC addresses, and 6104 unique IP addresses, a fingerprint database of 2078 standard TCP/IP signatures of 51 known unique OSes with a total of 529 variations when considering major and minor versions [30]. The dataset consists of three basic TCP/IP network stack features, i.e., initial SYN packet size, TTL, and TCP window size [30]. For the dataset with no label reductions, an accuracy of 84% was achieved. After our first set of testing, we realized that the data was severely skewed and that only a few of the classes contained almost all of the entries, giving us artificially good OSes classification results. In order to fix this issue, we explored two approaches: (1) keeping a bias, but with reduced differences, and (2) removing all bias and creating a dataset with equally distributed fingerprints for each class. For the first approach, the fewer occurring OS classes were copied until they reached the presence of at least half of the most occurring class. In the second approach, we removed most of the very seldom occurring classes and ended up with 33 reduced classes. We also removed all traffic that did not contain HTTP User-agent information, since we could not establish ground truth for this traffic. This led to the creation of a new fully general dataset where all the OS classes were bucketed into seven groups, consisting of the six most widely used major OS families: Android, Linux, Mac OS, Unix, Windows, iOS, and a seventh class called “Other” for OSes not suited for any of the other groups as shown in Table I. Finally, we ended up distributing all of the labels equally so that each OS class had the same number of occurrences. This is not necessarily the approach that gives a model with the best classification accuracy, but it creates the most versatile model with balanced training data. This helps us improve the generalizability of our model with a unified approach that encompasses all variations of the most widely used OSes.

TABLE II: The performance of OS classification from previous related work [30]

Method	Accuracy	Precision	Recall	F-score
User-agent	0.9189	0.9812	0.6063	0.7495
TCP/IP parameters	0.8088	0.5249	0.4643	0.4927
Specific domains	0.8402	0.6286	0.4907	0.5512
Combination	0.8582	0.6587	0.6041	0.6302

B. Realistic Traffic

While benchmark traffic is useful to link our experiments to previous related work, we also wanted additional realistic traffic for which we have more control, and that allows us to make our own assurances of the quality of the data. Thus, we passively collected our realistic dataset from TCP traffic originated from the internal network of the Oslo Metropolitan University and destined to various hosts on the Internet. First, we collected data for fixed (non-mobile) desktop computers (typically using OSes like Windows, Linux, Unix, Mac OSx, etc.) by using an intermediate node as shown in the network setup in Figure 1. Then, we passively collected the data that covered mobile devices, like *android* and *iOS*. The latter was collected from the 5G 4IoT research lab [1, 44] of the Oslo Metropolitan University.

We spent a significant amount of effort in establishing ground truth, i.e., determining the actual OS that has been used for each traffic flow. To establish ground truth in the realistic dataset, we follow two approaches. The first approach was only applicable to the non-mobile desktops, while the second method was used for both mobile and non-mobile devices. With the first method, we leveraged the DHCP log messages associated with the non-mobile desktops to derive the ground truth from the DHCP server of the Oslo Metropolitan University network that logs the sessions by the MAC address and name of the device. Since we collect the real data from the internal network of our university, extracting the DHCP log messages showing the distribution of MAC addresses and device names of the TCP sessions can give us detailed information about the OSes. The reason why we make use of these logs to determine the ground truth of the non-mobile devices is: since the network is dynamic, we cannot have full control over the connected devices. However, with a useful log of all the connected devices like this, we can make predictions and then compare the distribution of the predicted underlying OSes with the distribution of the MAC addresses and device names of the TCP sessions logged in the central DHCP server of our university.

Extracting the DHCP log messages can give us accurate matches and detailed information about the predicted underlying operating systems. We could, for example, see information about the *vendor-specific prefixes* since most of the OS variants are identified based on their vendors. The list of device vendor prefixes is useful in revealing the specific implementation of an OS because most of the modern OSes from the same device vendor usually share the same OS kernel and similar network behaviors. For example, we found out that Apple products often share the same TCP/IP parameters. This was the first approach we employed and it is so assuring though it takes a significant amount of effort. The second approach we used to identify the OS of the mobile and non-mobile devices is getting the predefined browser strings that loosely tell the name of the underlying OS assigned by the vendor from Webserver.

We believe changing the default device names by all users is not that common and sometimes discouraged by the vendors, e.g., Google and Apple OSes. However, the device name of Linux and Windows OSes could be changed easily by experienced users which would make passively identifying these devices hard. Since a number of computer vendors offer devices with a pre-installed OS and default device name and MAC address, we can use this information to derive the ground truth for OS fingerprinting. For example, Apple devices use a default string name of “<user>-iPhone”, “<user>-iPad”, Microsoft uses “Windows-Phone” for its mobile devices, and Android uses “android-<android_id>”, etc. Our real traffic covers the communication to and from our university and hence all traffic whose source and destination IP addresses are within the subnets of our internal network. Hence the network administrator of our university has full control over the internal machines with real IP addresses that are not going to a NAT gateway, and therefore it is fairly possible to tell whether it is a laptop or a desktop PC by looking it up in the internal database owned by the university. However, since

it is a dynamic network we do not have full control over external machines, because they can be anything behind an IP address that changes dynamically. This is because there is an endless number of machines spoofing scanning the network and they can appear as Linux-powered OSes but they could be Windows and vice versa and this happens because the user may have strongly tuned the TCP stack to look like something else. It is pretty hard to certainly say anything about the external computers because the communication can go through a NAT gateway possessing another OS type. For example, if a user is connected to a student wireless network, there is a chance that it may go to a Linux NAT gateway, and hence from outside the user is seen as Linux NAT which makes it hard to predict whether the underlying OS is Linux, Mac or Windows. Therefore, fingerprinting devices behind NAT technology on a distributed network where a number of devices can hide behind a NAT is another critical challenge. It is, therefore, worth noting that establishing ground truth in dynamic networks at a larger scale is a challenging problem and requires a lot of effort in the data preparation phase. Further investigation to explore these difficulties will be done in our future works. Finally, due to the privacy protection of possibly sensitive data, the payload of all the network packets collected was removed and anonymized with a prefix-preserving algorithm [10, 51]. Furthermore, we were only allowed to collect TCP headers of the traffic flows, while we could not collect complete traffic captures, due to privacy protection and legal reasons.

C. Emulated Traffic

In a real scenario where the OS fingerprinting is going on continuously in an intermediate node of an enterprise or production network, the intermediate node will have more information available than only the TCP header, such as the traffic profile or the knowledge of congestion or the outstanding bytes-in-flight of a TCP flow. In our experiments below, we show how this information can be very useful for OS fingerprinting. Since we do not have full traffic packet captures in our benchmark dataset or in our realistic dataset, we needed an additional dataset that we collected from an emulated network, where there would be no privacy protection or legal issues related to our dataset.

The architecture of our emulated network is similar to the network setup shown in Figure 1, except that all the nodes (the sender, the intermediate node, and the receiver) are implemented in virtual machines. All background traffic of the OSes for our emulated scenario is generated using the *iperf* [9]. In our experiment, the TCP flows were captured in the order of 1 minute duration and at a variable rate of 10 to 1000 Mbit/s. In order to build our model, we used about 3GB amount of training data. The setup to capture the traffic and the assumptions that we made are explained in further detail in our previous works [18, 19, 20]. Establishing ground truth in an emulated setup is straightforward, as we have full control of the OSes used when generating the traffic. In addition to establishing the ground truth, we also wanted to allow the intermediate node to establish a prediction of the TCP variant by monitoring the on-going traffic profile of the TCP flow

between the sender and the receiver. As shown later in the paper, using definitive or predicted knowledge of the TCP variant as an additional input feature to the OS fingerprinting, might boost the fingerprinting accuracy significantly. How the machine learning model for prediction of the TCP variant in the emulated scenario is trained and how the TCP variant is subsequently predicted are presented in the following section.

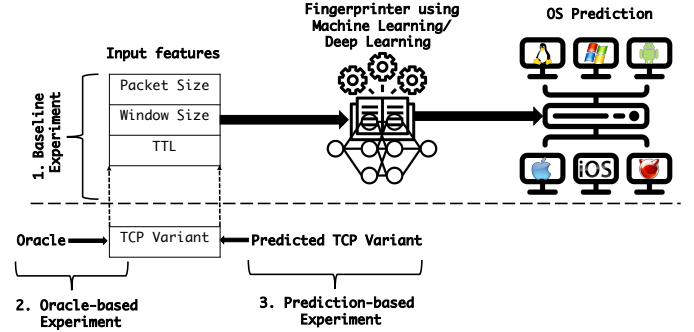


Fig. 2: The process implemented on the intermediate node for passive OS fingerprinting.

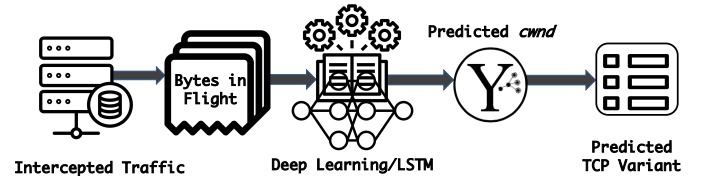


Fig. 3: The process implemented on the monitor for prediction of the TCP variant of the passively intercepted TCP traffic flow. An LSTM-based machine learning module predicts the *cwnd* from the outstanding bytes-in-flight. In the next step, the *cwnd* behavior is used to predict the underlying TCP variant as explained in further detail in our previous works [18, 19, 20]. As we can see in the bottom right part of Figure 2, The predicted TCP variant is finally used as an input feature to the OS fingerprinting process.

IV. MACHINE LEARNING OF THE OS FINGERPRINTER

A. Classical Machine Learning Approaches

The OS fingerprinter takes various features as input parameters, and use machine learning to predict the OS as shown in Figure 2. Many machine learning techniques could be used to implement a model for passive OS fingerprinting. In this paper, we have employed the following most commonly used classical machine learning methods suitable for our task. In order to train and test our classification models, we employed every experiment with a ratio of 60% training, 40% testing split, and 5-fold cross-validation setting on all variations of the features into one learning model.

SVM: In order to perform an efficient multi-class SVM classification through cross-validation, we tuned the SVM hyperparameters using a *GridSearchCV* that allows specifying

only the ranges of values for optimal parameters by parallelization construction of the model fitting. Finally, in our evaluation, we found out that *SVM* with a Radial Basis Function (RBF) kernel for classification model yields a substantially better result.

Random Forest (RF): We tuned the meta-estimator by varying the number of decision trees between 1 and 1000. We found out that increasing the number of trees more than 10 doesn't give much improvement in the classification accuracy.

KNN: We applied KNN by testing different values of K ranging from 5 to 100 followed by a weight function for a total of 20 observations. The observations have been conducted in two ways. In the first experiment, we set the weight to *uniform*. In the second experiment, the points are weighted by the inverse of their distance, causing closer neighbors to have greater influence. Finally, we choose the model that has the highest accuracy for a given unseen instance.

Naive Bayes (NB): Intuitively, advanced and modern machine learning methods are expected to perform better than classical techniques. Hence, in our experiment, we have employed classical classification machine learning methods like NB model as a baseline classifier. As it is shown in the experimental results, given its simplicity and effectiveness, it consistently performs comparably well as the other traditional classification models with a small inaccuracy margin.

B. Deep Learning Approaches

To find the deeper characteristics of TCP variants implemented by respective OSes and exploit the extra OS-specific information, we apply the following two neural network architectures.

Multilayer Perceptron (MLP). In our evaluation, MLP model with a single-layer feedforward neural network [22, 42] has been used to classify the different classes of OSes. After the hyperparameter tuning, we tested our MLP model with a different number of batch sizes, hidden layers, and nodes (e.g., 0, 1, 2, 32, 64, 128) in each layer. Combining all of these, a total of 324 models were trained with and without the default TCP variant. We found out that the results for both with and without a known TCP variant were almost the same with an insignificant drop in the accuracy irrespective of which hyperparameters performed the best. Finally, 150 nodes of the network per dataset are trained for 500 epochs with a batch size of 500 by SGD with momentum of 0.9 and a constant learning rate of 0.01. However, we learned that SGD is sensitive in regards to the selection of the learning rate since it doesn't automatize the values and we also found that it suffers from premature convergence and is outperformed by *Adam*-based optimization methods. Hence, both *Adam* and *Nadam* gradient-based optimization algorithms fit for our purpose and that is because we wanted to use an optimization algorithm that adapts its learning rate dynamically in a way that doesn't affect the objective function and learning process of the model. Our experimental results show that the hyperparameter tuning baseline experiments by applying *tanh* as activation

function and *Adam* optimization algorithm and training the model for 500 epochs, provides a substantial improvement in accuracy as compared to the other parameters.

Long Short-Term Memory (LSTM) models. We have explored an approach to classify the underlying OS from passive measurements using LSTM-based RNN architecture by combining the basic TCP/IP features and the underlying TCP variant shown in Table 2 as input vectors. For more details about LSTM applied in the context of computer networks, we refer the reader to our previous paper [19]. We trained our LSTM model over 500 epochs of the training samples with a batch size of 250 as values in time-series. We propagate the input feature vector (x) to the model through a multilayer LSTM cell followed by a fully connected dense layer of 150 hidden nodes with Rectified Linear Unit (ReLU) activation function using the *hard_sigmoid* as recurrent activation for the different layers that generates an output of a sequence dimensional vector of predicted OSes (y_t). We trained our LSTM-based learning algorithm without the knowledge of the input features from the true signatures of the OSes during the learning phase. We learn the model from the training data and then finally predict the test labels from the testing instances on all variations of the OS-specific parameters. In order to get a more stable and robust to changes of the passive OS classification model, we have applied the *Adam* stochastic optimizer algorithm [26]. It is one of the most effective optimization algorithms in the deep learning community. In our experiment, the algorithm is set with an initial *learning rate* of 0.001 and *exponential decay rates* of the first (β_1) and second (β_2) moments set to 0.9 and 0.999 respectively. We further optimize a wide range of important hyperparameters related to the neural network topology to improve the performance of our passive OS classification model.

C. Comparative Suitability

Here, we demonstrate the comparative suitability of implementing each classical machine learning and deep learning classifiers analyzed for the benchmark, emulated, and realistic datasets we used in our paper.

SVM is often used as a baseline technique for both binary and multi-class classification tasks in the machine learning community. In addition to this, *SVM* classifiers use *kernelization* in order to handle non-linearly separable features. *SVM* techniques also work fine with unbalanced datasets like the benchmark data we used in our paper. In our experiment, we employed kernel *SVM* with RBF for classification equipped with different kernel functions and regularization parameters. RF classification models, as compared to *SVM* classifiers, have fewer problems handling non-linearity. Furthermore, RF models perform slightly better when it comes to high-dimensional regression and classification tasks. Another main reason why we used RF as an evaluation approach in our experiments is that it is relatively fast and it is possible to maintain a reasonably acceptable accuracy with inconsistent and unbalanced datasets. Moreover, state-of-the-art KNN models for neighborhood-based classification are very simple,

effective, and also handle both linearly and non-linearity separable features reasonably well. Bayesian methods, on the other hand, are incapable of handling well non-linearly separable features. NB operates with the strong hypothesis of statistical independence between the features of the model. However, the reason why we have employed NB in our evaluation as a baseline classifier is to compare the performance of modern and older machine learning methods. Advanced deep learning models have great potential for handling non-linearity in the dataset but at the cost of longer training time as can be seen in Table XVII and by introducing a heavy computational burden.

D. Experimental Hardware Setup

All our machine learning experiments are carried out using a cluster of HPC machines based upon the GNU/Linux operating system running a modified version of the 4.15.0-39-generic kernel release. The prediction model is performed on an NVIDIA Tesla K80 GPU accelerator computing with the following characteristics: Intel(R) Xeon(R) CPU E5-2670 v3 @2.30GHz, 64 CPU processors, 128 GB RAM, 12 CPU cores running under Linux 64-bit. All nodes in the cluster are connected to a low latency 56 Gbit/s Infiniband, gigabit Ethernet, and have access to 600 TiB of BeeGFS parallel file system storage.

E. Objectives of Our Experiments

The aim of our experiments is to explore the role of the underlying TCP variant as an input feature when passively detecting the underlying OS. To investigate this, we divide our analysis into three different experiments.

First, in the baseline experiment presented in Section V, we carry out the OS fingerprinting without using a known TCP variant as an input feature. This corresponds to the simplest state-of-the-art transport layer method, which is illustrated in the upper part of Figure 2. Since there is a close connection between existing popular OSes and the TCP variants they use, our hypothesis was that the potential for improvement by using the TCP variant as an input feature would be significant. For example, CUBIC [15] is the default congestion control algorithm as part of the Linux kernel distribution configurations from version 2.6.19 onwards. Since Android devices are also Linux-powered, CUBIC remains to be the default TCP congestion control algorithm. Many Windows 7 distributions have been shipped with the default New Reno [21] and whereas Windows 8 families with CTCP [48]. Therefore, in the next Oracle-based experiment presented in Section V, we investigate the potential of knowing the TCP variant, and how much this knowledge might boost the fingerprinting accuracy. Here we assume that there is an Oracle that can identify and give the TCP variant used in the TCP flow that is fingerprinted. This is illustrated in the bottom left part of Figure 2. However, in a real scenario, the intermediate node would not have access to definite knowledge of the TCP variant (e.g., given by an Oracle). Instead, the intermediate node might at best try to infer it from the monitored traffic. Thus, in the third prediction-based experiment presented in Section VI, we

first allow the intermediate node to predict the TCP variant passively. This is illustrated in the bottom right part of Figure 2. The OS fingerprinter then uses that TCP variant prediction as an input feature to make the OS prediction illustrated in the upper part of Figure 2. The TCP variant is predicted by analyzing the famous sawtooth pattern behavior of estimated cwnd of TCP, which is computed based on the outstanding bytes-in-flight [19, 20]. This is presented in more detail in the next section. Since the latter experiment requires TCP traffic details of outstanding bytes-in-flight, which is not available in our benchmark and realistic datasets, this experiment is only possible with our emulated dataset.

V. THE ROLE OF TCP VARIANT IN PASSIVE OS FINGERPRINTING AND ITS FEASIBILITY

In this section, we perform primary experiments on whether the underlying TCP variant is a relevant input feature in passive OS fingerprinting.

A. Baseline experiment: Results without knowing the TCP variant

Here we present the results of the machine learning and deep learning techniques under all the validation scenarios presented above without a known underlying TCP variant which will play the role of baseline for the other evaluations.

Based on benchmark data from previous related work:

Looking at Table III, both machine learning and deep learning classification techniques have consistently achieved good levels of precision and recall for all general classes of OSes except iOS. Quantitatively, iOS, and Mac OS devices were underrepresented in the benchmark data from previous related work. Besides, as it is shown in Figure 4, there is a slightly higher misclassification of iOS as unknown and this is why the precision and recall of iOS are comparably lower than the rest of OSes. We also believe that the limited TCP/IP stack basic features could contribute to the indistinguishability and misclassification of OS classes with the same kernel implementation. The false positives are easier to notice in the corresponding confusion matrices. As discussed above, since the benchmark dataset presented in [30] is skewed we argue that using balanced accuracy is a better indicator of performance for such datasets as shown in Tables III and VI.

Based on realistic traffic: Our performance results of the realistic traffic without a known TCP variant using the machine learning and deep classification techniques are presented in Table IV. The respective normalized confusion matrix for each technique are presented in Figure 5.

Based on emulated traffic: Our performance results of the emulated traffic without a known TCP variant as an input feature using both machine learning and deep learning techniques are presented in Table V. As we can see in the corresponding normalized confusion matrices presented in Figure 6, the precision and recall for most of the OSes using both machine learning and deep learning techniques are reasonably good.

TABLE III: Benchmark data [30] experimental results without a known TCP variant using machine learning and deep learning OS classification techniques.

OS	Machine Learning Techniques								Deep Learning Techniques			
	SVM		RF		KNN		NB		MLP		LSTM	
	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall
Android	0.75	0.89	0.86	0.90	0.86	0.90	0.68	0.85	0.75	0.92	0.77	0.85
Linux	0.84	0.83	0.91	0.91	0.86	0.94	0.83	0.82	0.90	0.82	0.83	0.85
Mac OS	0.64	0.76	0.61	0.82	0.59	0.83	0.64	0.76	0.62	0.81	0.62	0.83
Other	0.91	0.81	0.91	0.81	0.91	0.81	0.88	0.81	1.00	0.74	0.91	0.81
Unix	0.94	0.99	0.94	0.99	0.94	0.99	0.92	0.99	0.94	0.99	0.94	0.99
Windows	0.97	0.89	0.98	0.89	0.98	0.89	0.98	0.79	0.97	0.91	0.97	0.86
iOS	0.71	0.54	0.71	0.54	0.79	0.48	0.69	0.54	0.67	0.57	0.79	0.55
<i>Average</i>	0.82	0.82	0.85	0.84	0.85	0.83	0.80	0.79	0.84	0.82	0.83	0.82
Balanced Accuracy	82.17%		84.76%		84.28%		80.11%		82.98%		81.87%	

TABLE IV: Realistic traffic experimental results without a known TCP variant using machine learning and deep learning OS classification techniques.

OS	Machine Learning Techniques								Deep Learning Techniques			
	SVM		RF		KNN		NB		MLP		LSTM	
	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall
Android	0.75	0.89	0.86	0.90	0.84	0.93	0.68	0.85	0.81	0.83	0.76	0.86
Linux	0.89	0.82	0.94	0.89	0.93	0.88	0.85	0.82	0.89	0.79	0.90	0.81
Mac OS	0.63	0.81	0.61	0.82	0.61	0.82	0.64	0.76	0.61	0.82	0.82	0.79
Unix	0.94	0.99	0.94	0.99	0.94	0.99	0.92	0.99	0.92	0.99	0.94	0.99
Windows	0.97	0.89	0.98	0.89	0.98	0.89	0.98	0.82	0.98	0.89	0.97	0.89
iOS	0.88	0.72	0.86	0.73	0.88	0.72	0.86	0.72	0.84	0.73	0.70	0.92
<i>Average</i>	0.85	0.83	0.86	0.85	0.87	0.85	0.83	0.81	0.84	0.83	0.83	0.84
Accuracy	83.43%		85%		85.10%		81.25%		83.91%		83.27%	

TABLE V: Emulated traffic experimental results without a known TCP variant using machine learning and deep learning OS classification techniques.

OS	Machine Learning Techniques								Deep Learning Techniques			
	SVM		RF		KNN		NB		MLP		LSTM	
	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall
Android	0.74	0.90	0.86	0.90	0.85	0.91	0.74	0.88	0.75	0.88	0.91	0.85
Linux	0.92	0.82	0.94	0.89	0.92	0.90	0.84	0.85	0.93	0.78	0.92	0.74
Mac OS	0.63	0.81	0.61	0.82	0.61	0.82	0.64	0.76	0.62	0.81	0.86	0.88
Unix	0.94	0.99	0.94	0.99	0.94	0.99	0.94	0.99	0.92	0.99	0.94	1.00
Windows	0.97	0.89	0.98	0.89	0.98	0.89	0.97	0.88	0.93	0.91	0.98	0.73
iOS	0.88	0.73	0.86	0.73	0.88	0.73	0.88	0.73	0.88	0.73	0.82	1.00
<i>Average</i>	0.85	0.84	0.86	0.85	0.87	0.85	0.84	0.83	0.85	0.83	0.89	0.88
Accuracy	84.67%		85.73%		85.27%		83.12%		84.05%		88.44%	

TABLE VI: Benchmark data [30] experimental results with Oracle-given TCP variant using machine learning and deep learning OS classification techniques.

OS	Machine Learning Techniques								Deep Learning Techniques			
	SVM		RF		KNN		NB		MLP		LSTM	
	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall
Android	0.96	1.00	0.99	0.98	0.99	0.98	0.93	0.98	0.98	0.96	0.96	0.99
Linux	0.86	0.93	0.92	0.95	0.91	0.95	0.82	0.92	0.87	0.94	0.90	0.93
Mac OS	0.99	0.90	0.96	0.92	0.96	0.92	0.98	0.88	0.96	0.92	0.99	0.90
Other	0.93	0.81	0.93	0.81	0.91	0.83	0.91	0.81	0.93	0.81	1.00	0.74
Unix	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
Windows	0.99	0.89	0.97	0.91	0.99	0.89	1.00	0.78	0.97	0.91	0.83	0.92
iOS	0.75	0.88	0.75	0.91	0.75	0.91	0.71	0.89	0.76	0.89	0.72	0.85
<i>Average</i>	0.92	0.92	0.93	0.93	0.93	0.92	0.91	0.89	0.92	0.92	0.91	0.90
Balanced Accuracy	91.83%		92.56%		92.33%		90.05%		91.97%		91.58%	

Comparison of results without known TCP variant: As shown in Tables III, IV, and V, our experimental results are pretty consistent. Firstly, we can see that there is not much difference in performance across different machine learning and deep learning techniques. But more importantly, there are not many differences in performance between results from

using different types of experimental data. This is intuitively correct, since the OS fingerprinting is based on the basic TCP/IP packet fields, and should not differ much between various types of data, whether we do evaluation using the benchmark data, real data or emulated data. Secondly, we believe accuracy in the range of 82-88% (average value) is

TABLE VII: Realistic traffic experimental results with Oracle-given TCP variant using machine learning and deep learning OS classification techniques.

OS	Machine Learning Techniques								Deep Learning Techniques			
	SVM		RF		KNN		NB		MLP		LSTM	
	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall
Android	0.95	1.00	0.99	0.98	0.99	0.98	0.93	1.00	0.97	1.00	0.97	0.97
Linux	0.86	0.91	0.94	0.93	0.92	0.94	0.83	0.91	0.91	0.92	0.90	0.93
Mac OS	0.99	0.90	0.96	0.92	0.97	0.92	1.00	0.88	0.99	0.90	0.97	0.90
Unix	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
Windows	0.99	0.89	0.99	0.89	0.99	0.89	1.00	0.84	0.99	0.89	0.99	0.89
iOS	0.93	0.96	0.91	0.99	0.92	0.98	0.91	0.96	0.91	0.98	0.92	0.97
Average	0.95	0.95	0.96	0.96	0.96	0.96	0.94	0.94	0.95	0.95	0.96	0.95
Accuracy	94.81%		95.65%		95.69%		93.62%		95.12%		95.14%	

TABLE VIII: Emulated traffic experimental results with the Oracle-given TCP variant using machine learning and deep learning OS classification techniques.

OS	Machine Learning Techniques								Deep Learning Techniques			
	SVM		RF		KNN		NB		MLP		LSTM	
	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall
Android	0.97	0.98	0.99	0.98	0.99	0.98	0.95	1.00	0.98	0.97	0.96	0.98
Linux	0.90	0.91	0.95	0.93	0.92	0.95	0.88	0.90	0.97	0.89	0.93	0.91
Mac OS	0.99	0.90	0.97	0.92	0.97	0.92	0.99	0.90	0.93	0.94	0.94	0.92
Unix	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
Windows	0.99	0.89	0.97	0.91	0.97	0.91	0.99	0.89	0.99	0.89	0.98	0.88
iOS	0.91	0.98	0.92	0.98	0.93	0.97	0.91	0.97	0.91	0.99	0.91	0.97
Average	0.95	0.95	0.96	0.96	0.96	0.96	0.95	0.95	0.95	0.95	0.95	0.95
Accuracy	95.10%		96.02%		95.83%		94.60%		95.24%		95.08%	

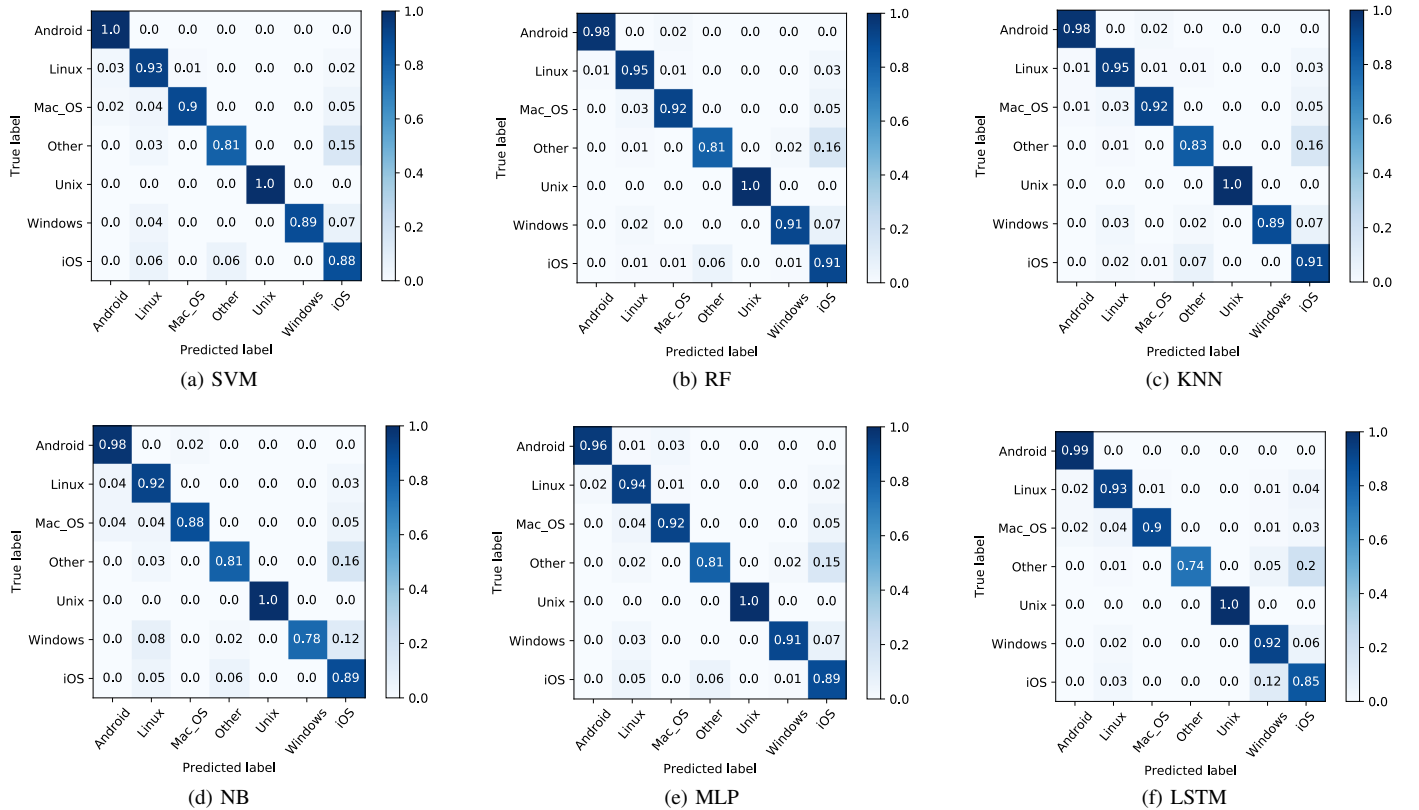


Fig. 4: Normalized confusion matrix comparison of the machine learning and deep learning OS classification techniques with Oracle-given TCP variant using the benchmark data from previous related work [30].

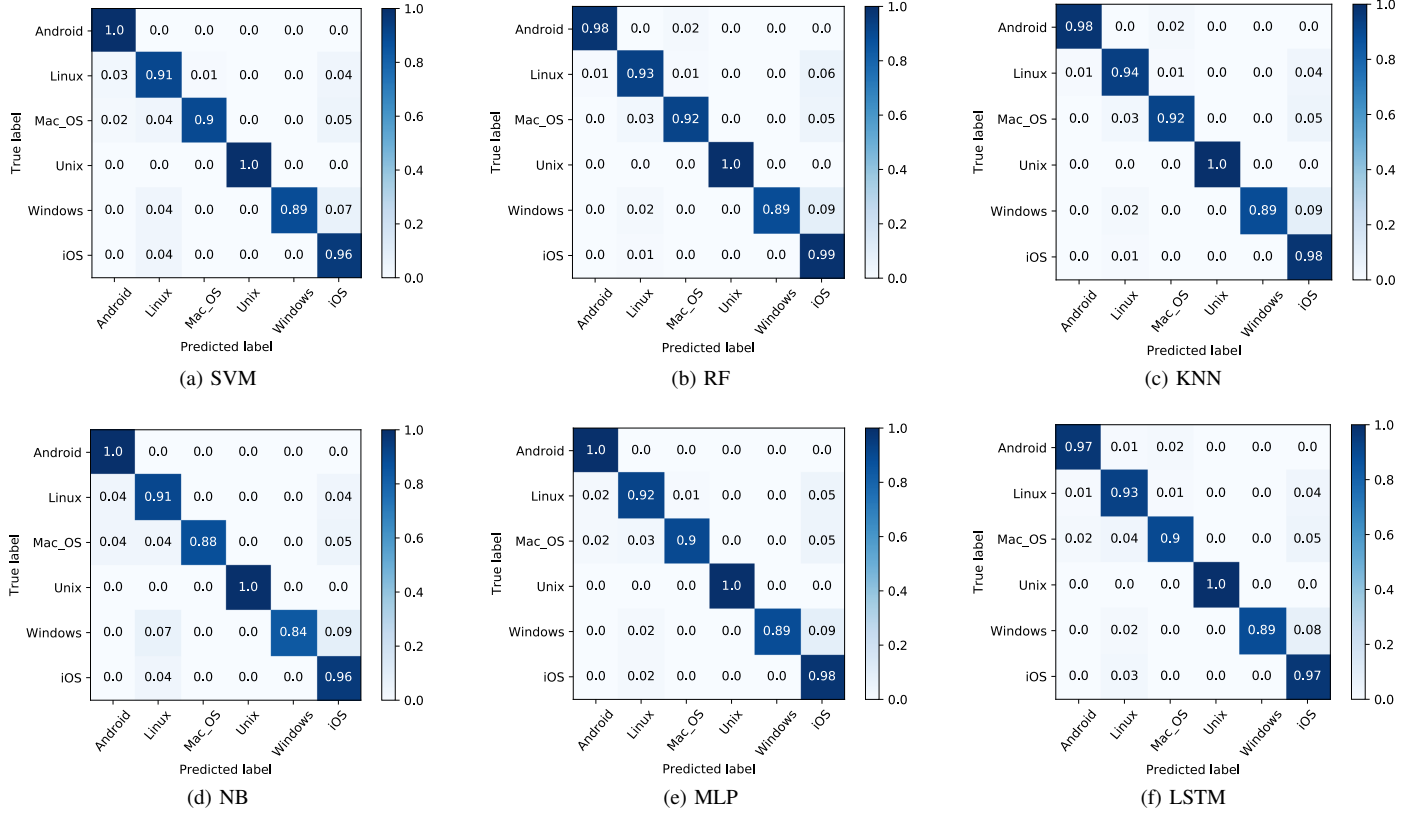


Fig. 5: Normalized confusion matrix comparison of the machine learning and deep learning OS classification techniques with Oracle-given TCP variant using a realistic traffic.

perhaps not sufficient for a product in a real deployment. Our hypothesis is that this accuracy could be boosted considerably had we only known the implementation of the underlying TCP variant. We will explore this hypothesis in the next section.

B. Oracle-based experiment: Results using Oracle-given TCP variant

Here we assume that we know exactly the underlying TCP variant, i.e., we assume it is given by an Oracle. We show that knowledge of the TCP variant has a great potential for boosting passive fingerprinting of OSes, and in this section, we will try to quantify this potential. In the next section, we will show that much of this potential can be harvested by using a tool that predicts the underlying TCP variant.

Based on benchmark data from previous related work: Table VI shows a significant performance gain across all classes of OSes when we assume prior knowledge of the underlying TCP variant, as compared to the results when the TCP variant is unknown presented in Table III.

Based on realistic traffic: The performance results of the realistic traffic with the Oracle-given TCP variant presented in Table VII shows the potential of knowing TCP variant given by an Oracle for passive OS fingerprinting in a realistic scenario.

Based on emulated traffic: Our performance results of the emulated traffic with the Oracle-given TCP variant using both classical machine learning and deep learning techniques are presented in Table VIII. We can see that this shows a significant improvement in performance over the results without a known TCP variant presented in Table V. Both machine learning and deep learning techniques have comparable and consistent results in terms of accuracy.

Comparison of results with Oracle-given TCP variant: Our accuracy results presented in Tables III, VII, and VIII, demonstrate that by knowing the TCP variant we obtain a considerable performance boost in all our experimental results, compared to our previous results obtained without knowledge of the TCP flavor. With an Oracle-given TCP variant, we obtain a prediction accuracy of 94-96%, with an average value of 94.1% over all traffic classes and of 95.4% over only emulated traffic. The accuracy results are pretty consistent across all scenarios. Comparing these results with our previous results that do not use the Oracle (84.1% on average for all traffic types and 85.6% only for emulated traffic), we observe a solid increase in the OS fingerprinting performance. This improvement would significantly boost the usefulness of a product to be implemented in a real enterprise network infrastructure. As in the previous section, here again, we observe highly consistent performance results across different

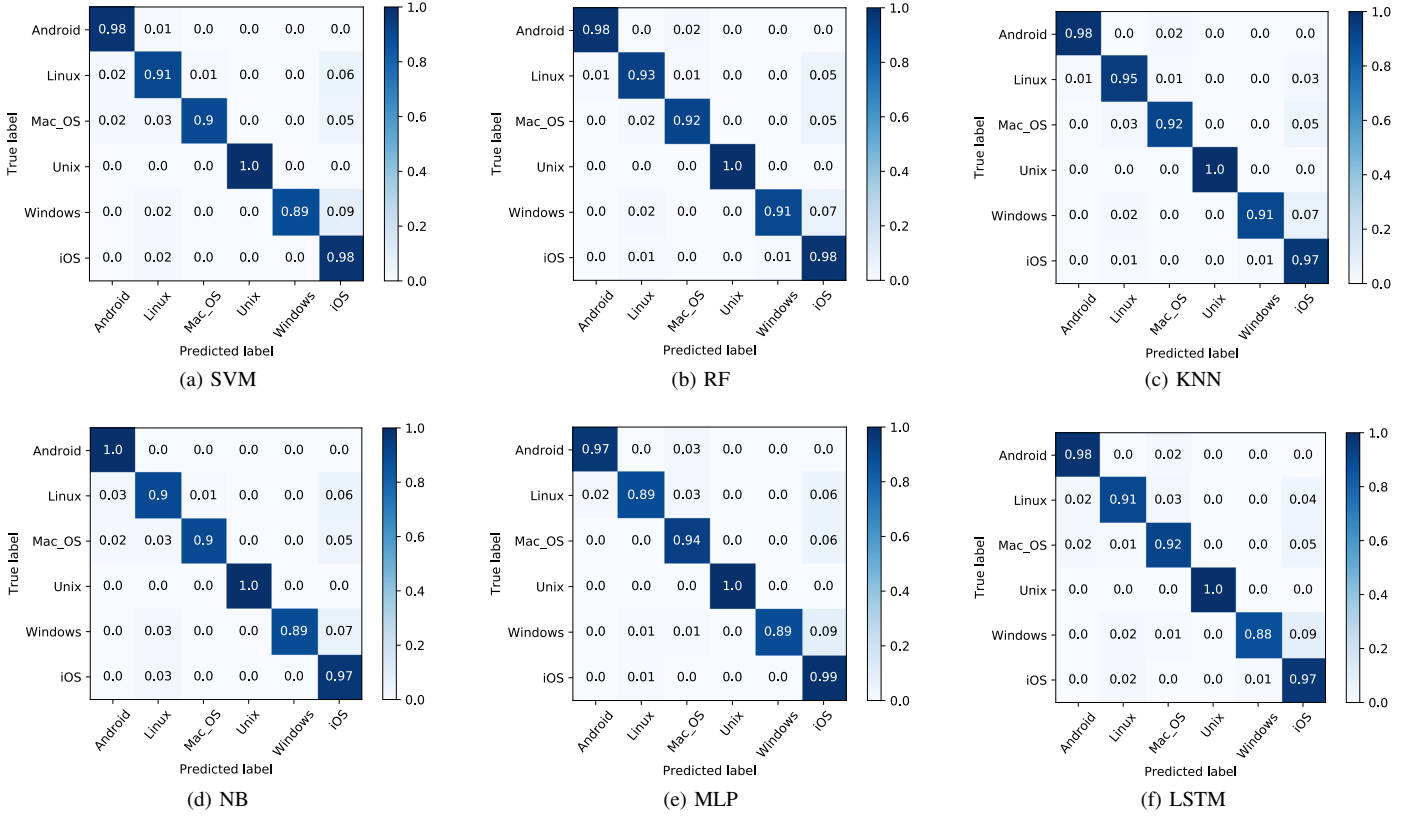


Fig. 6: Normalized confusion matrix comparison of the machine learning and deep learning OS classification techniques with Oracle-given TCP variant using emulated traffic.

machine learning and deep learning techniques and also between the use of different types of experimental data. The latter is useful knowledge for the next section since it means that performance increases obtained over one traffic type is shown to be amenable to other traffic types as well.

In the next section, we will have to base our evaluation on emulated data, since we do not have the TCP traffic patterns of the realistic data or benchmark data at hand. These traffic patterns are required to be able to passively infer the TCP variant in the experiments presented in the next section. In this section, the idealistic Oracle was used only to demonstrate the potential of knowing the TCP variant, but this is not a realistic assumption. Thus, in the next section, we will instead base our evaluation on a TCP variant that is passively predicted by a deep learning-based tool that we developed and presented in our previous work [18, 19, 20]. Using this tool, we explore how close our performance will get to the ideal solution of having an Oracle-given TCP variant.

Feature Selection: For the traditional machine learning algorithms, we can employ any feature selection algorithm that can be used to experiment with the performance of each input feature. Feature selection is an NP-hard problem and the vast majority of those feature selection algorithms employ some greedy criteria to select a subset of features. However, since we have very few input features (4 features in our case),

we opted for a more systematic and computationally feasible approach where we checked all possible combinations of 2 features and 3 features to yield an optimal feature selection as shown in Tables IX and X. This could help us understand the impact of each input feature on the passive OS classification performance. Interestingly, the TCP flavor was present in the combination of 2 features and 3 features and it consistently improves the classification performance as shown Tables IX and X. We could employ this approach for all the classical machine learning algorithms presented in our paper. However, to avoid redundancy for the reader, the feature combinations presented in Tables IX and X are only for the RF algorithm. Effective passive OS fingerprinting analysis requires more variations in network traffic. As it is specified in [41] and [40], inspecting a combination of the TTL in the IP header and the size of the TCP receiver window of the first packet in a TCP session is often enough in order to successfully fingerprint various OSes of target remote computers. One main reason behind why the values of TTL and TCP receiver window size vary is that different OSes and different versions of the same underlying OS set different default values for these parameters [40, 41]. As it can be seen from the experimental results presented in Tables IX and X, a combination of the input features TTL and the TCP window size achieves a better accuracy consistently when we have both two and three pairs of the input features.

TABLE IX: RF classification results of emulated traffic using two features combinations of the initial SYN packet size (PS), TCP receiver window size (WS), TTL, and the Oracle-given TCP variant.

OS	Feature Combinations											
	PS, WS		PS, TTL		PS, Oracle-given TCP		WS, TTL		WS, Oracle-given TCP		TTL, Oracle-given TCP	
	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall
Android	0.72	0.77	0.78	0.59	0.68	0.90	0.72	0.89	0.79	0.88	0.66	0.90
Linux	0.91	0.75	0.80	0.84	0.91	0.74	0.86	0.82	0.90	0.84	0.89	0.75
Mac OS	0.62	0.78	0.60	0.78	0.64	0.76	0.64	0.76	0.61	0.82	0.64	0.76
Unix	0.94	0.99	0.94	0.99	0.94	0.99	0.94	0.99	0.94	0.99	0.94	0.99
Windows	0.94	0.77	0.92	0.80	0.95	0.86	0.95	0.87	0.95	0.88	0.94	0.81
iOS	0.75	0.74	0.75	0.74	0.85	0.73	0.88	0.73	0.88	0.73	0.85	0.73
<i>Average</i>	0.80	0.79	0.79	0.78	0.83	0.81	0.84	0.83	0.85	0.84	0.83	0.81
Accuracy	79.8%		78.2%		81.4%		82.6%		83.7%		81%	

TABLE X: RF classification results of emulated traffic using three features combinations of the initial SYN packet size (PS), TCP receiver window size (WS), TTL, and the Oracle-given TCP variant.

OS	Feature Combinations							
	PS, WS, TTL		PS, TTL, Oracle-given TCP		WS, TTL, Oracle-given TCP		PS, WS, Oracle-given TCP	
	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall
Android	0.86	0.90	0.92	0.91	0.97	0.98	0.92	0.97
Linux	0.94	0.89	0.94	0.82	0.89	0.91	0.92	0.87
Mac OS	0.61	0.82	0.91	0.88	0.97	0.92	0.97	0.88
Unix	0.94	0.99	1.00	1.00	1.00	1.00	1.00	1.00
Windows	0.98	0.89	0.87	0.86	0.94	0.89	0.98	0.83
iOS	0.86	0.73	0.88	0.96	0.92	0.96	0.87	0.98
<i>Average</i>	0.86	0.85	0.92	0.91	0.95	0.94	0.93	0.93
Accuracy	85.73%		91.4%		94.5%		92.9%	

C. TCP variant prediction tool

The main goal of the experiments in the emulated network is to use the predicted TCP variant as an additional distinguishing input feature to the passive OS fingerprinting. The TCP variant is predicted by the process illustrated in Figure 3. As described in sufficient detail in our previous works [18, 19, 20], we used a database to match and join the intercepted TCP traffic on both the intermediate node and the sending node. The outstanding bytes-in-flight of the traffic (i.e., the number of bytes that have been sent but not yet acknowledged) is used as input to our machine learning model to predict the cwnd behaviour of the traffic. We use LSTM for the machine learning. We trained and verified the machine learning model by matching the predicted TCP states with the actual TCP kernel states directly logged from the Linux kernel of the sending node. Since we have full control of the sending nodes, we can track the system-wide TCP state of every packet that is sent and received from the kernel to verify our model’s prediction accuracy against the actual TCP variant by matching with the actual sending TCP states using the techniques presented in our previous works [18, 19, 20].

After the verification, we can run our learning model and get the cwnd predictions of the TCP stack in use. Once we can estimate the cwnd of the sender, we can also infer the multiplicative back-off factor to decrease the cwnd on a loss event (β) which is an important feature for uniquely identifying the TCP variants. Finally, we combine the predicted TCP variant as the basis of OS fingerprinting with the basic TCP/IP features as shown in Figure 2. Here, we consider only loss-based TCP congestion control algorithms, e.g., BIC [52], CUBIC [15], CTCP [48], Reno [23], and New Reno [21]. Our approach could also be useful to other TCP variants like

Google’s QUIC [29]. QUIC is a general-purpose transport layer network protocol that uses packet loss as an indicator of congestion and supports a number of different congestion control algorithms, including CUBIC [15] and BBR [5].

VI. PASSIVE OS FINGERPRINTING BASED ON PREDICTED TCP VARIANT

Based on the types of implicit congestion signals and other local information, the underlying TCP congestion control algorithms are categorized into loss-based and delay-based variants. In this section, we demonstrate the potential of passively predicting both the loss-based and delay-based TCP flavors for improving the passive OS fingerprinting.

A. Results using loss-based TCP variant prediction

In Section V, we showed that Oracle-given knowledge of the underlying TCP variant has a great potential for improving the passive OS fingerprinting. In reality, however, we don’t have an Oracle-given TCP variant. Since passively detecting the TCP variant is a challenging task, this is where our tool from previous works on predicting the underlying TCP variant from passive measurements [18, 19, 20] comes into play. In this Section we use the TCP variant passively *predicted* by this tool as an input feature for the passive OS fingerprinting. The TCP variant is inferred from the famous Additive Increase and Multiplicative Decrease (AIMD) sawtooth pattern of TCP’s estimated cwnd computed based on the outstanding bytes-in-flight. Since we don’t have access to the actual cwnd of the senders in the benchmark data and realistic traffic, here we consider only the emulated traffic.

Based on emulated traffic: In this section, we use a tool to predict the TCP variant from passive measurements of TCP traffic patterns, and this prediction is used as input to the OS fingerprinting method presented above. The experimental results of both techniques are presented in Table XI.

Comparison of results with a predicted TCP variant: Results with emulated data and a passive prediction of the TCP variant as shown in Table XI reveal an accuracy of 91.2% on average, which comes pretty close to the accuracy of 95.3% obtained on emulated traffic with the TCP-variant given by the Oracle. Intuitively, when we perform learning based on the TCP variant prediction, the OS classification accuracy must be lower than the Oracle-given TCP variant. But the question is how close can we get to the idealistic scenario of having an Oracle. Our results show that using our tool for TCP variant prediction from passive measurements gives reasonably good OS fingerprinting accuracies that come close to the results obtained by using the Oracle-given TCP variant. Even though the performance results with the TCP variant passively predicted by our deep learning-based tool are slightly lower as compared to the TCP variant given by an idealistic Oracle, our performance results of using our tool are reasonably competitive.

B. Results using delay-based TCP protocols

The passive OS fingerprinting method presented above, where the cwnd is first computed based on the outstanding bytes-in-flight, then the underlying TCP flavor is predicted from the estimated cwnd, is particularly efficient for loss-based TCP variants that consider packet loss as an implicit indication of congestion. Unlike traditional loss-based TCP variants, delay-based TCP congestion control algorithms use the changes in queueing delay measurements as implicit feedback to congestion in the network. Delay-based congestion control algorithms attempt to avoid network congestion by monitoring the trend of network path's Round-Trip Time (RTT) information contained in packets [16]. By design, unlike loss-based TCP algorithms, the multiplicative decrease parameter (β) of delay-based congestion control algorithms is not fixed which makes it fundamentally challenging to predict the TCP variant from passive traffic measurements when there is variability in delay. For example, TCP VenO [13] sets β factor to 0.8 when the queueing delay is small. However, when the queueing delay is high, TCP VenO [13] sets β to 0.5. The back-off parameter along with other TCP characteristics can be used to predict the underlying TCP congestion control algorithms. In our previous work [17], we have developed an efficient tool for the prediction of the underlying delay-based TCP flavors from passive measurements by utilizing the β and queueing delay values. By using different data-driven classification techniques based on probabilistic models and Bayesian inference approaches, we addressed how the β varies as a function of queueing delay changes and investigated into how the TCP variants of delay-based congestion control algorithms can be predicted both from passively measured traffic and real measurements over the Internet [17].

In this section we will extend the passive OS fingerprinting method presented above by coupling to our previous work [17] to also cover delay-based TCP variants, e.g., TCP Vegas [4], TCP VenO [13], BBR [5], etc. The performance results with emulated data and a passive prediction of the delay-based TCP flavors as presented in Table XII show an accuracy of 95.24% and 95.38% on average using both classical machine learning and deep learning techniques respectively. The corresponding confusion matrices of these techniques are presented in Figure 7. We can see that this shows a significant improvement in performance over the results without a known TCP variant presented above in Table VIII. Both machine learning and deep learning techniques have comparable and consistent results in terms of accuracy. Our experimental results show that using our statistical methods for delay-based TCP flavors prediction gives reasonably good OS fingerprinting accuracies. However, in a realistic assumption, we don't know exactly if the sending node is using either a loss-based or delay-based TCP flavor implementation. Therefore, we need a generic passive OS fingerprinting tool that can take both loss-based and delay-based TCP variants as input, and make a reasonably good OS classification.

To make our passive OS fingerprinting tool generic, we run a separate extensive experiment in an emulated setting with a combination of loss-based (e.g., TCP Reno [23], BIC [52], and CUBIC [15]) and delay-based (e.g., TCP VenO [13] and TCP Vegas [4]) TCP variants. As a result, we obtain an OS fingerprinting performance accuracy of 94.95% and 95.22% on average using machine learning and deep learning techniques respectively as shown in Table XIII. This shows that our model can also be applied equally well to loss-based as well to delay-based TCP variants as input. As in the previous sections, here again, we observe both our machine learning and deep learning classification techniques under an emulated setting have consistently achieved good levels of precision and recall for all general classes of OSes. By combining different variants of TCP, we demonstrate that our passive OS fingerprinting tool is generic enough which gives promising and comparable results in terms of accuracy across different experimental scenarios.

C. Transfer Learning Results

One of the primary benefits of employing machine learning and deep learning techniques as discussed in Section I is the concept of *transfer learning*. In the machine learning community, *transfer learning* is defined as the ability to take a model trained in one experiment scenario and apply it for classification in a different experiment scenario. For example, in our case, that means we are able to train our model on a dataset created in an emulated network with an Oracle-given TCP variant and apply it for classification of our dataset from the realistic traffic. Results presented in Table XIV shows that the learning of the OS fingerprinter using loss-based TCP variants transfers well into other scenarios. Similarly, as it can be seen from the results shown in Table XV, the learning of our OSes fingerprinting model using both loss-based and delay-based predicted TCP variants transfers

TABLE XI: Emulated traffic experimental results with loss-based predicted TCP variant using machine learning and deep learning OS classification techniques.

OS	Machine Learning Techniques								Deep Learning Techniques			
	SVM		RF		KNN		NB		MLP		LSTM	
	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall
Android	0.92	0.96	0.92	0.97	1.00	0.97	0.92	0.95	0.95	0.97	0.92	0.96
Linux	0.79	0.85	0.94	0.82	0.92	0.94	0.80	0.89	0.98	0.79	0.86	0.90
Mac OS	0.96	0.88	0.97	0.87	0.85	0.94	0.97	0.88	0.95	0.90	0.95	0.88
Unix	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
Windows	0.92	0.78	0.85	0.80	0.88	0.91	0.98	0.65	0.94	0.77	0.97	0.77
iOS	0.85	0.94	0.86	0.96	0.93	0.87	0.84	0.95	0.82	0.99	0.88	0.96
<i>Average</i>	0.90	0.90	0.91	0.91	0.93	0.93	0.91	0.90	0.92	0.91	0.92	0.92
Accuracy	90.01%		91.09%		92.15%		90.40%		91.45%		91.93%	

TABLE XII: Emulated traffic experimental results with delay-based predicted TCP variant using machine learning and deep learning OS classification techniques.

OS	Machine Learning Techniques								Deep Learning Techniques			
	SVM		RF		KNN		NB		MLP		LSTM	
	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall
Android	0.95	1.00	0.98	0.98	0.99	0.98	0.95	1.00	0.97	0.97	0.98	0.95
Linux	0.88	0.90	0.91	0.94	0.94	0.93	0.86	0.91	0.96	0.88	0.93	0.91
Mac OS	0.99	0.90	0.98	0.92	0.98	0.92	0.99	0.90	0.93	0.94	0.95	0.92
Unix	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.99	1.00
Windows	0.99	0.89	0.98	0.88	0.99	0.89	0.99	0.89	0.97	0.91	0.98	0.89
iOS	0.91	0.97	0.92	0.98	0.91	0.99	0.93	0.96	0.91	0.97	0.90	0.98
<i>Average</i>	0.95	0.95	0.95	0.95	0.96	0.96	0.95	0.95	0.95	0.95	0.95	0.95
Accuracy	94.60%		95.86%		95.81%		94.68%		95.61%		95.14%	

TABLE XIII: Emulated traffic experimental results with a combination of loss-based and delay-based predicted TCP variant using machine learning and deep learning OS classification techniques.

OS	Machine Learning Techniques								Deep Learning Techniques			
	SVM		RF		KNN		NB		MLP		LSTM	
	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall
Android	0.97	0.98	0.97	0.98	0.99	0.98	0.94	0.99	0.98	0.98	0.98	0.97
Linux	0.91	0.91	0.91	0.90	0.91	0.95	0.89	0.90	0.93	0.91	0.96	0.87
Mac OS	0.99	0.90	0.98	0.92	0.97	0.92	0.97	0.91	0.94	0.94	0.92	0.94
Unix	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.99	1.00
Windows	0.99	0.89	0.98	0.90	0.99	0.89	0.98	0.88	0.99	0.89	0.95	0.90
iOS	0.91	0.98	0.91	0.97	0.92	0.98	0.90	0.96	0.92	0.98	0.91	0.97
<i>Average</i>	0.95	0.95	0.95	0.95	0.96	0.96	0.94	0.94	0.95	0.95	0.95	0.94
Accuracy	95.04%		95.28%		95.78%		93.70%		95.37%		95.07%	

TABLE XIV: Transfer learning experimental results with loss-based predicted TCP variant using machine learning and deep learning OS classification techniques.

OS	Machine Learning Techniques								Deep Learning Techniques			
	SVM		RF		KNN		NB		MLP		LSTM	
	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall
Android	0.95	1.00	0.98	0.98	0.99	0.98	0.94	0.99	0.97	0.98	0.97	0.96
Linux	0.86	0.91	0.90	0.95	0.92	0.95	0.85	0.93	0.95	0.85	0.91	0.91
Mac OS	0.99	0.90	0.98	0.92	0.97	0.92	0.99	0.90	0.94	0.94	0.96	0.90
Unix	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
Windows	0.99	0.89	0.98	0.90	0.97	0.91	0.99	0.84	0.99	0.89	0.98	0.87
iOS	0.93	0.96	0.93	0.97	0.93	0.97	0.90	0.95	0.90	0.98	0.90	0.98
<i>Average</i>	0.95	0.95	0.95	0.95	0.96	0.96	0.94	0.93	0.95	0.95	0.94	0.94
Accuracy	94.79%		95.35%		95.76%		93.54%		94.72%		94.28%	

well across other scenarios. A transfer learning experiment combining the loss-based and delay-based predicted TCP variants for an OS fingerprinting as presented in Table XV gives an accuracy of 94.83% and 94.88% on average using both classical machine learning and deep learning techniques respectively. The corresponding normalized confusion matrix for both machine learning and deep learning techniques is shown in Figure 8. Good transfer learning results indicate that

our passive OS fingerprinting model is able to discern the results of unforeseen scenarios and still perform reasonably well. In our previous works, we have also demonstrated that the TCP variant predictor performs well in terms of transfer learning [18, 19, 20]. In summary, this shows that our multi-class classification model is general bearing similarity to the concept of transfer learning in the machine learning community [37, 38, 50].

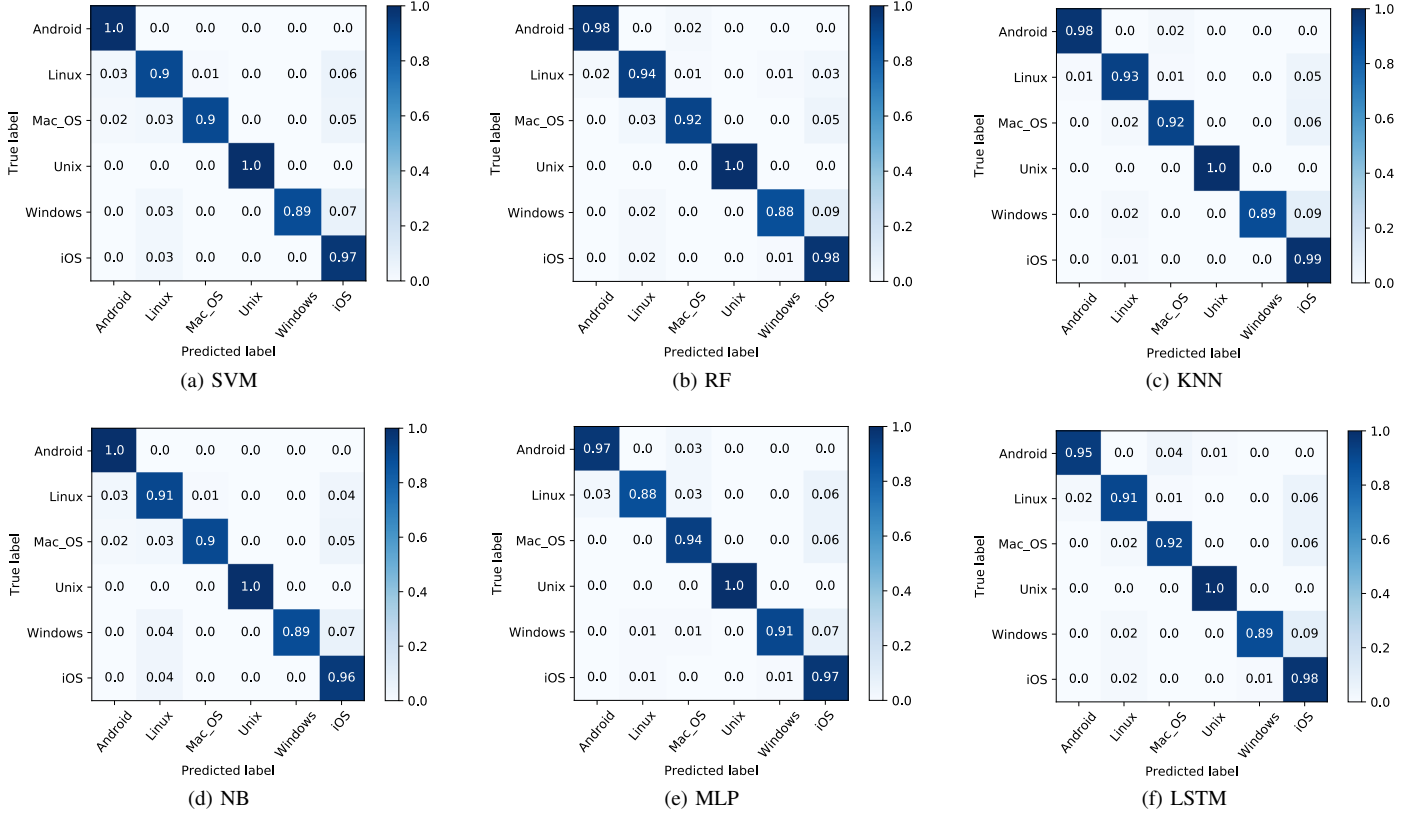


Fig. 7: Normalized confusion matrix comparison of the classical machine learning and deep learning OS classification techniques for predicted delay-based TCP variants using emulated traffic.

TABLE XV: Transfer learning experimental results with a combination of loss-based and delay-based predicted TCP variant using machine learning and deep learning OS classification techniques.

OS	Machine Learning Techniques								Deep Learning Techniques			
	SVM		RF		KNN		NB		MLP		LSTM	
	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall
Android	0.96	1.00	0.96	1.00	0.99	0.98	0.92	0.97	0.98	0.95	0.98	0.97
Linux	0.88	0.90	0.91	0.91	0.94	0.94	0.89	0.91	0.94	0.88	0.94	0.90
Mac OS	0.99	0.90	0.99	0.90	0.97	0.92	0.97	0.88	0.93	0.94	0.93	0.94
Unix	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.99	1.00
Windows	0.99	0.99	0.97	0.90	0.99	0.89	0.94	0.89	0.97	0.91	0.99	0.87
iOS	0.91	0.89	0.92	0.97	0.91	0.98	0.92	0.96	0.91	0.97	0.91	0.98
Average	0.95	0.97	0.95	0.95	0.96	0.96	0.94	0.94	0.95	0.95	0.95	0.95
Accuracy	94.83%		94.99%		95.74%		93.78%		94.90%		94.86%	

D. Discussion

Comparison of our approach with other OS fingerprinting tools: Here, the accuracy of our machine learning and deep learning approaches presented in this paper are compared against the *state-of-the-art* passive OS fingerprinting tool, p0f [56]. Table XVI presents the comparison of our approaches and p0f evaluated using emulated traffic under different settings. As we can see from Table XVI, even though the performances are reasonably comparable, the experimental results show that our passive OS fingerprinting approaches outperform the *state-of-the-art* p0f method across all scenarios except without a known TCP variant.

TABLE XVI: Performance comparison of our approaches with p0f using emulated traffic under different settings.

	Loss-based TCP variants		
	Our approach		P0f
	Machine Learning	Deep Learning	
Without a known TCP	84.70%	86.25%	88.12%
Oracle-given TCP	95.39%	95.16%	90.03%
Predicted TCP	90.91%	91.69%	90.89%
Delay-based TCP variants			
Predicted TCP	95.24%	95.38%	91.73%
Loss-based and delay-based TCP variants			
Predicted TCP	94.95%	95.22%	91.38%

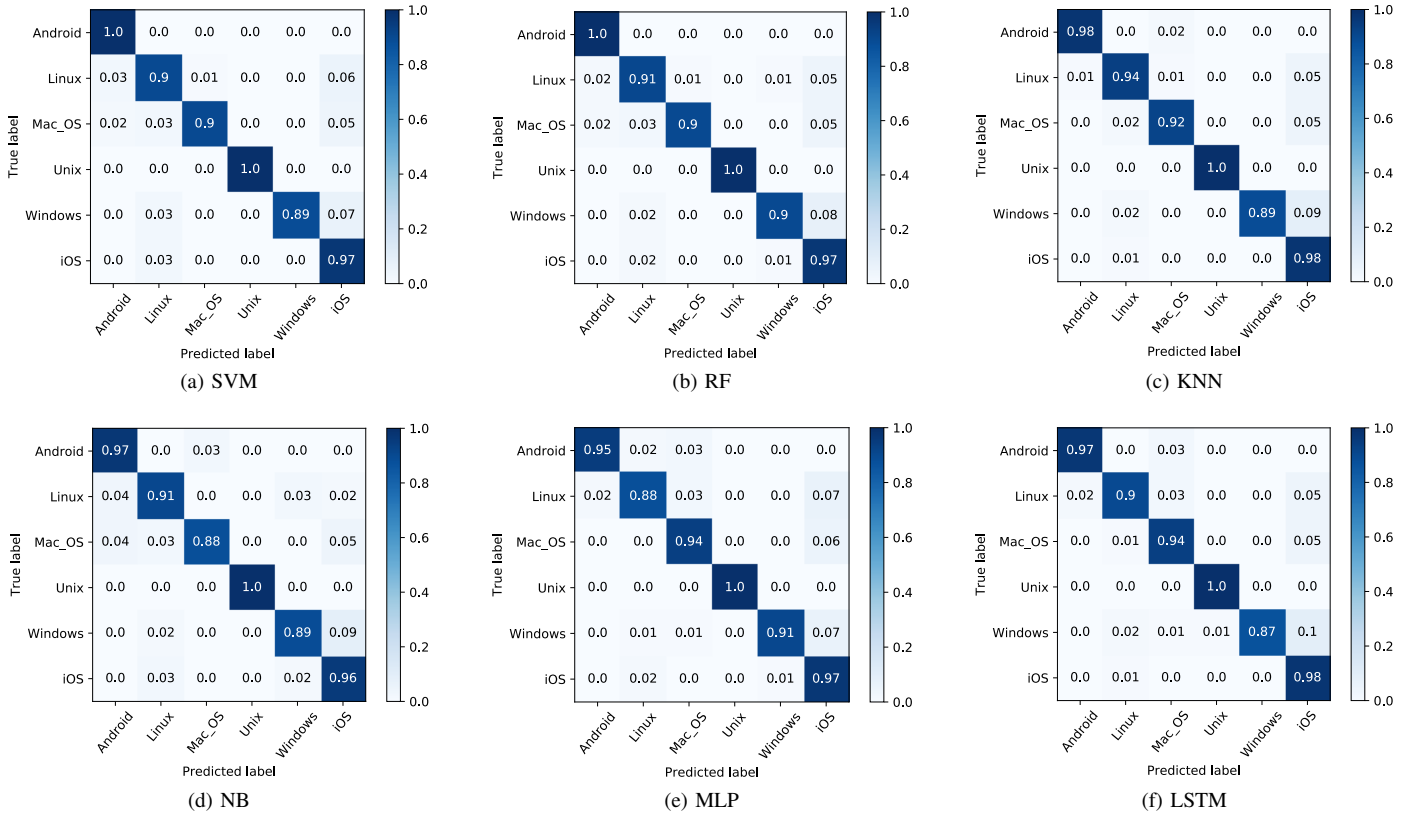


Fig. 8: Transfer learning: Normalized confusion matrix comparison of the classical machine learning and deep learning OS classification techniques for loss-based and delay-based predicted TCP variants.

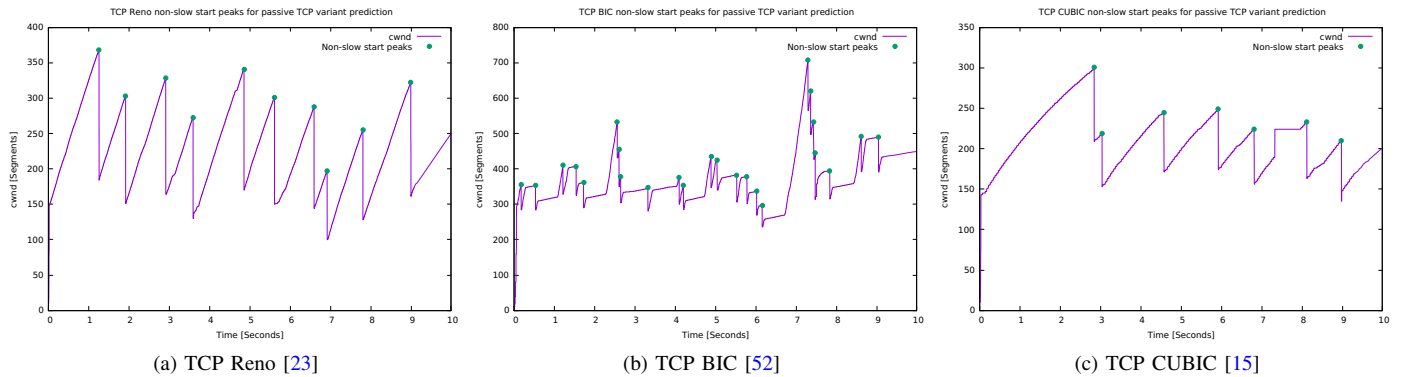


Fig. 9: Non-slow start peaks for passive TCP variant predictions of (a) Reno (b) BIC and (c) CUBIC using emulated traffic.

Execution time complexity: The complexity of the approaches presented in our paper in terms of how long it takes (in seconds) to passively perform OS fingerprinting is presented in Table XVII. The execution time of each machine learning and deep learning techniques for all the scenarios we considered in our paper is using the same NVIDIA Tesla GPU requirements. The execution time for the SVM model is significantly longer than the rest of traditional machine

learning classification techniques and this is because the computing kernel values of SVM takes more time and memory requirements. Irrespective of their execution time, the ability in dealing with unbalanced data is one of the reasons why we included both SVM, RF, and KNN classification methods in our experiment. NB, as compared to the other machine learning and deep learning models, takes relatively a small amount of time to train the prediction model as it can be

TABLE XVII: Execution time (seconds) of the machine learning and deep learning techniques.

		Loss-based TCP variants					
		SVM	RF	KNN	NB	MLP	LSTM
Benchmark	Without a known TCP variant	2258.2201247215	1.1143145561	1.9390485286	1.1081347370	112.7487759590	1806.6641292572
	Oracle-given TCP variant	1120.1224200725	1.1178588867	2.0935020446	1.1160198802	114.5473096370	1777.935036697
Realistic	Without a known TCP variant	8789.4385774135	7.2072355747	9.9768719673	6.1865331881	573.7393236160	5753.8345386981
	Oracle-given TCP variant	3726.3648753166	7.2179739475	11.4665050506	7.1971922931	577.4705853462	5702.6515867710
Emulated	Without a known TCP variant	3937.6824579238	4.1011998653	6.9211487770	4.1011255760	250.6548788547	748.2380180358
	Oracle-given TCP variant	1571.1400921344	4.1095368862	5.1454861164	3.1089884179	251.8693726062	652.0555405616
	Predicted TCP variant	975.8222060203	3.9419437313	5.9578386783	3.1154929704	19.2053611278	673.2168228626
Delay-based predicted TCP variants							
Emulated	Predicted TCP variant	3819.2830920219	4.3223826885	5.7427815914	3.2945567131	484.2410390377	1305.2965178489
Loss-based and delay-based predicted TCP variants							
Emulated	Predicted TCP variant	11348.1196880340	4.2940688133	9.2223489284	4.2719212722	733.8166005611	3743.1047005653
Transfer Learning							
Loss-based TCP Variants		935.2456374168	5.1156289577	6.0417027473	4.1121685368	250.9682075977	2973.5155694484
Loss-based and delay-based TCP variants		8255.4148607254	5.3129029273	9.9110872745	4.3107805252	786.6981770992	1980.7358047962

TABLE XVIII: TCP variant prediction accuracies on emulated scenarios with a different number of peaks.

	The first one peak		The first two peaks		The first three peaks		The first five peaks		All peaks (the whole flow)	
	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall
BIC	1.00	0.53	1.00	0.65	1.00	0.85	0.97	0.97	0.97	1.00
CUBIC	0.70	0.97	0.73	0.97	0.81	0.97	0.92	0.97	0.97	1.00
Reno	0.57	0.61	0.66	0.66	0.83	0.79	0.94	0.89	1.00	0.95
<i>Average</i>	0.75	0.70	0.79	0.76	0.88	0.87	0.94	0.94	0.98	0.98
Accuracy	70.37%		75.93%		87.04%		94.44%		98.15%	

TABLE XIX: TCP variant prediction confusion matrices on emulated scenarios with a different number of peaks.

Actual	Predicted														
	The first one peak			The first two peaks			The first three peaks			The first five peaks			All peaks (the whole flow)		
	BIC	CUBIC	Reno	BIC	CUBIC	Reno	BIC	CUBIC	Reno	BIC	CUBIC	Reno	BIC	CUBIC	Reno
BIC	18	0	16	22	0	12	29	0	5	33	0	1	35	0	0
CUBIC	0	35	1	0	35	1	0	35	1	0	35	1	0	35	0
Reno	0	15	23	0	13	25	0	8	30	1	3	34	1	1	36

seen in Table XVII since it is faster to train. As we can see from all our experimental results, the performance of MLP and LSTM models is reasonably comparable. However, the execution time of MLP is relatively smaller than the LSTM model as seen in Table XVII. The reason why the LSTM model takes much longer average execution time is that MLP is much more efficient.

Amount of traffic for TCP variant prediction: In order to determine the amount of traffic our approach requires to passively infer the underlying loss-based and delay-based TCP variants of clients, we consider the whole flow of the TCP session since it significantly improves the overall prediction performance. As shown in Table XVIII, we carried out an experiment under a different number of peaks to test how our approach infers the loss-based TCP variant using emulated traffic. We based our peak analysis for predicting the TCP variant on the multiplicative back-off parameter, β , value by averaging out the window size of AIMD algorithm every time we have a peak so that we don't do the computation of the multiplicative decrease factor only on a slow-start phase. There are two approaches to measure the β value of a TCP congestion control algorithm: (i) when there is a packet loss event, and (ii) when a time out event occurs. The β value, along with other TCP characteristics, especially for loss-based congestion control algorithms is one of the most important parameters

which determines important conditions of network congestion like the cwnd and slow-start threshold (sssthresh) [53]. Hence, as shown in Figure 9, we use the β value so as to uniquely predict the underlying TCP variant based on the multiplicative back-off factor of the selected loss-based TCP variants. According to the TCP standard specification, the β value of Reno [23], BIC [52], and CUBIC [15] is set to fixed values of 0.5, 0.8, and 0.7 respectively. The TCP variant prediction accuracies on emulated scenarios with the first one, two, three, five, all the peaks of the flow are 70.37%, 75.93%, 87.04%, 94.44%, and 98.15% respectively as shown in Table XVIII and their corresponding confusion matrix is depicted in Table XIX. As we can see from these results, it is clear that the higher the number of peaks we consider for the analysis, the better the TCP variant prediction is and this the reason why we argue considering the whole flow of the TCP session is a better approach. Unlike loss-based algorithms, it is worth noting that delay-based TCP congestion control algorithms, by design, have a variable β and the β value of these protocols varies when there is variability in queuing delay which makes it fundamentally challenging to predict the TCP variant from passive traffic. In our previous work [17], we have presented an effective TCP variant identification methodology that addresses how β varies as a function of queuing delay and how delay-based TCP variants can be predicted both from passively measured traffic and real measurements over the Internet.

E. Limitations and possible improvements of our approach

Here we detail some limitations of our approach and possible improvements to address them. In our method, we have used the flow duration, which is in the order of 1 minute in our experiments (see Section III for more details), as the granularity in the experiment. However, in practice, our approach is rather dependent on observing enough number of TCP events between endpoints in order to accurately recognize the underlying TCP variant.

The current trend on the Internet is to use multiple streams carrying HTTP traffic between the same endpoints to avoid head of line blocking. This implies that the duration of flow will be reduced in the future. However, the total number of packets between endpoints will not change and only the duration will be reduced. Hence, the number of TCP events on multiple flows between the same endpoints, not the duration of a single flow, is significant for fingerprinting, which is worth investigating in future experiments. We believe that our proposed approach can be easily adapted to accommodate TCP variant identification after observing enough number of congestion events between endpoints instead of relying on a single TCP flow which might be too short to observe those events. Furthermore, in this paper, we only consider the potential of the underlying TCP variant as a distinguishable input feature, but for future work, we believe that a more comprehensive TCP-based feature which includes the variations of the TCP cwnd computed in one stage can form in itself a unique signature of an OS since the standard implementations even of the same TCP congestion control mechanisms differ [11].

VII. CONCLUSION AND FUTURE WORK

Passively fingerprinting the underlying OS implementation of a remote host is important for security-conscious network administrators. It can, for example, be used in identifying the source of malicious traffic, exploring a network for potential exploitations of security vulnerabilities, defining OS-based access control security policies, configuring network-based IDS to classify and prioritize extraneous security alerts etc. In this paper, we proposed and evaluated a novel approach that attempts to passively fingerprint the underlying remote OS by leveraging *state-of-the-art* machine learning and deep learning classification techniques under multiple controlled scenarios. We show that knowing the Oracle-given TCP variant has a great potential for boosting the classification performance of passive OS fingerprinting. In our setting, we demonstrate that using the idealistic Oracle has the potential to boost the prediction accuracy from 84.1% to 94.1% on average across all traffic types tested, and from 85.6% to 95.3% in an emulated setting. However, in reality, we don't have the Oracle-given TCP variant and hence we don't know exactly the underlying TCP flavor. To address this, we demonstrated a method for passive OS fingerprinting where the cwnd is first computed based on the outstanding bytes-in-flight, then the underlying TCP flavor is predicted from the estimated cwnd, and finally, the predicted TCP variant is used as an input feature to detect the remote computer's OS. This is an

additional feature that is added to the basic TCP/IP features that are the basis of OS fingerprinting in previous works. We demonstrate that our method performs significantly better than not using the predicted underlying TCP variant as an input feature, increasing the accuracy in our experiment from 85.6% to 91.2% and 95.3% on average using loss-based and delay-based TCP variants respectively.

By combining both loss-based and delay-based predicted TCP flavors, our OS fingerprinting model achieves an accuracy of 95.22%. The results of this method come close to the accuracy of 95.4% obtained by using the idealistic Oracle. To the best of our knowledge, this is the first study that reports the potential of the underlying TCP feature in boosting significantly the accuracy of passive OS fingerprinting. We further validate and demonstrate the transferability approach of our OSes classification models by conducting a series of controlled experiments against other scenarios. Through comparing the experimental results between the benchmark dataset, realistic, and emulated traffic in terms of accuracy and confusion matrix, it is clear that our passive OSes classification models are able to discern the results to unforeseen scenarios. Therefore, we are able to show that the learned passive OS fingerprinting model by leveraging a pre-trained knowledge of classification techniques from the emulated network performs reasonably well as it is shown in the experimental results when it is applied and transferred to a realistic scenario. Lastly, in all our experiments, we made sure that both the training and validation accuracies are closer which gives an idea about the ability of the OSes classification models to generalize on unforeseen scenarios.

Note that passively detecting the underlying TCP variant is fundamentally a challenging task, which led to a two-step approach in our paper, where the TCP variant prediction of a deep learning-based tool is used as input to another machine learning method in the next step. However, by integrating the two machine learning approaches better, there should be potential for increasing the passive OS classification performance even further and get even closer to the idealistic results of using an Oracle-given TCP variant. Exploring such optimizations is also left for future work. It is known that TCP clock drift improves OS fingerprinting and hence measuring differences in the timing of how the IP stack works may allow us to predict the underlying OS with greater assurance in terms of accuracy. We, therefore, argue for using other TCP options like timestamps and queueing delay characteristics as an input feature vector for passive OSes fingerprinting model as another interesting direction.

Finally, in addition to the difficulties of establishing ground truth (e.g., the underlying TCP variant) at a larger scale on a dynamic network addressed in Section III, there is a lot of other work to be done as an extension of our work presented here. As a future work, since our proposed method relies on passively identifying the underlying TCP variant accurately, we aim to study the vulnerability of our method to adversarial changes in the behavior of the underlying TCP flavor which usually requires a lot of user expertise.

ACKNOWLEDGMENT

We would like to thank the Norwegian center for research data (NSD) for granting us the legal permission to collect a realistic experiment dataset that contains OS-specific information from the Oslo Metropolitan University network. We would also like to thank the 5G 4IoT research lab at Oslo Metropolitan University for allowing us to collect the realistic data for android and iOS devices.

REFERENCES

- [1] 5G4IoT. 5G4IoT. <http://5g4iot.vlab.cs.hioa.no/>, 2019.
- [2] A. Aksoy, S. Louis, and M. H. Gunes. Operating system fingerprinting via automated network traffic analysis. In *IEEE Congress on Evolutionary Computation (CEC)*. IEEE, 2017.
- [3] I. Andrea, C. Chrysostomou, and G. Hadjichristofi. Internet of Things: Security vulnerabilities and challenges. In *2015 IEEE Symposium on Computers and Communication (ISCC)*, pages 180–187. IEEE, 2015.
- [4] L. S. Brakmo, S. W. O’Malley, and L. L. Peterson. *TCP Vegas: New techniques for congestion detection and avoidance*, volume 24. ACM, 1994.
- [5] N. Cardwell, Y. Cheng, C. S. Gunn, S. H. Yeganeh, and V. Jacobson. BBR: Congestion-based congestion control. *Queue*, 14(5):20–53, 2016.
- [6] D. Chang, Q. Zhang, and X. Li. Study on OS Fingerprinting and NAT/Tethering based on DNS Log Analysis. In *IRTF & ISOC Workshop on Research and Applications of Internet Measurements (RAIM)*, 2015.
- [7] W. R. Cheswick, S. M. Bellovin, and A. D. Rubin. *Firewalls and Internet security: repelling the wily hacker*. Addison-Wesley Longman Publishing Co., Inc., 2003.
- [8] N. Davids. Initial TTL values. http://noahdavids.org/self_published/TTL_values.html, 2011.
- [9] ESnet. iperf3. <https://iperf.fr/iperf-servers.php>, 2017.
- [10] J. Fan, J. Xu, M. H. Ammar, and S. B. Moon. Prefix-preserving IP address anonymization: measurement-based security evaluation and a new cryptography-based scheme. *Computer Networks*, 46(2):253–272, 2004.
- [11] P. Fiterău-Broștean, R. Janssen, and F. Vaandrager. Combining model learning and model checking to analyze TCP implementations. In *International Conference on Computer Aided Verification*, pages 454–471. Springer, 2016.
- [12] J. Franklin, D. McCoy, P. Tabriz, V. Neagoie, J. V. Randwyk, and D. Sicker. Passive Data Link Layer 802.11 Wireless Device Driver Fingerprinting. In *USENIX Security Symposium*, volume 3, pages 16–89, 2006.
- [13] C. P. Fu and S. C. Liew. TCP Venio: TCP enhancement for transmission over wireless access networks. *IEEE Journal on selected areas in communications*, 21(2):216–228, 2003.
- [14] L. G. Greenwald and T. J. Thomas. Toward Undetected Operating System Fingerprinting. *WOOT*, 7:1–10, 2007.
- [15] S. Ha, I. Rhee, and L. Xu. CUBIC: a new TCP-friendly high-speed TCP variant. *ACM SIGOPS operating systems review*, 42(5):64–74, 2008.
- [16] D. H. Hagos, P. E. Engelstad, and A. Yazidi. A Deep Learning Approach to Dynamic Passive RTT Prediction Model for TCP. In *IEEE 38th International Performance Computing and Communications Conference (IPCCC)*. IEEE, 2019.
- [17] D. H. Hagos, P. E. Engelstad, and A. Yazidi. Classification of Delay-based TCP Algorithms From Passive Traffic Measurements. In *2019 IEEE 18th International Symposium on Network Computing and Applications (NCA)*, pages 1–10. IEEE, 2019.
- [18] D. H. Hagos, P. E. Engelstad, A. Yazidi, and Ø. Kure. A machine learning approach to TCP state monitoring from passive measurements. In *2018 Wireless Days (WD)*, pages 164–171. IEEE, 2018.
- [19] D. H. Hagos, P. E. Engelstad, A. Yazidi, and Ø. Kure. Recurrent Neural Network-based Prediction of TCP Transmission States from Passive Measurements. In *2018 IEEE 17th International Symposium on Network Computing and Applications (NCA)*, pages 1–10. IEEE, 2018.
- [20] D. H. Hagos, P. E. Engelstad, A. Yazidi, and O. Kure. Towards a Robust and Scalable TCP Flavors Prediction Model from Passive Traffic. In *2018 27th International Conference on Computer Communication and Networks (ICCCN)*, pages 1–11. IEEE, 2018.
- [21] T. Henderson, S. Floyd, A. Gurtov, and Y. Nishida. The NewReno modification to TCP’s fast recovery algorithm. RFC 6582, 2012.
- [22] K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5):359–366, 1989.
- [23] V. Jacobson. Congestion avoidance and control. In *ACM SIGCOMM computer communication review*. ACM, 1988.
- [24] V. Jacobson, R. Braden, and D. Borman. TCP extensions for high performance. RFC 1323, 1992.
- [25] Q. Jing, A. V. Vasilakos, J. Wan, J. Lu, and D. Qiu. Security of the Internet of Things: perspectives and challenges. *Wireless Networks*, 20(8):2481–2501, 2014.
- [26] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [27] T. Kohno, A. Broido, and K. C. Claffy. Remote physical device fingerprinting. *IEEE Transactions on Dependable and Secure Computing*, 2(2):93–108, 2005.
- [28] E. Kollmann. Chatter on the Wire: A look at DHCP traffic. [Online]. Available: <http://myweb.cableone.net/xnih/download/chatter-dhcp.pdf> [Accessed: May 19, 2010], 2007.
- [29] A. Langley, A. Riddoch, A. Wilk, A. Vicente, C. Krasic, D. Zhang, F. Yang, F. Kouranov, I. Swett, J. Iyengar, et al. The quic transport protocol: Design and internet-scale deployment. In *Proceedings of the Conference of the ACM Special Interest Group on Data Communication*, pages 183–196. ACM, 2017.
- [30] M. Lastovicka, T. Jirsik, P. Celeda, S. Spacek, and D. Filakovsky. Passive os fingerprinting methods in the jungle of wireless networks. In *NOMS 2018-2018 IEEE/IFIP Network Operations and Management Symposium*, pages 1–9. IEEE, 2018.
- [31] R. Lippmann, D. Fried, K. Piwowarski, and W. Streilein. Passive operating system identification from TCP/IP packet headers. In *Workshop on Data Mining for Computer Security*, volume 40. Citeseer, 2003.
- [32] R. Lippmann, S. Webster, and D. Stetson. The effect of identifying vulnerabilities and patching software on the utility of network intrusion detection. In *International Workshop on Recent Advances in Intrusion Detection*, pages 307–326. Springer, 2002.
- [33] G. F. Lyon. Remote OS detection via TCP/IP stack fingerprinting. *Phrack Magazine*, 8(54), 1998.
- [34] G. F. Lyon. *Nmap network scanning: The official Nmap project guide to network discovery and security scanning*. 2009.
- [35] Netresec. Networkminer. <https://www.netresec.com/?page=NetworkMiner>, 2007.

- [36] A. Ornaghi and M. Valleri. Ettercap. <https://www.ettercap-project.org/>, 2015.
- [37] S. J. Pan. Transfer Learning., 2014.
- [38] S. J. Pan and Q. Yang. A survey on transfer learning. *IEEE Transactions*, 2010.
- [39] J. Postel. Internet control message protocol. RFC 792, 1981.
- [40] J. Postel. Internet protocol. RFC 791, 1981.
- [41] J. Postel. Transmission control protocol. RFC 793, 1981.
- [42] F. Rosenbaltt. The perceptron—a perceiving and recognizing automation. *Technical Report 85-460-1 Cornell Aeronautical Laboratory*, 1957.
- [43] J. Scambray, S. McClure, and G. Kurtz. *Hacking exposed*. McGraw-Hill Professional, 2000.
- [44] SCOTT. European Leadership Joint Undertaking. <https://scottproject.eu/>, 2019.
- [45] S. Sicari, A. Rizzardi, L. A. Grieco, and A. Coen-Porisini. Security, privacy and trust in Internet of Things: The road ahead. *Computer networks*, 76:146–164, 2015.
- [46] R. Spangler. Analysis of remote active operating system fingerprinting tools. *University of Wisconsin*, 2003.
- [47] G. Taleck. Synscan: Towards complete tcp/ip fingerprinting. *CanSecWest, Vancouver BC, Canada*, pages 1–12, 2004.
- [48] K. Tan, J. Song, Q. Zhang, and M. Sridharan. A compound TCP approach for high-speed and long distance networks. In *Proceedings IEEE INFOCOM 2006. 25TH IEEE International Conference on Computer Communications*, pages 1–12. IEEE, 2006.
- [49] W. Wei, K. Suh, B. Wang, Y. Gu, J. Kurose, and D. Towsley. Passive online rogue access point detection using sequential hypothesis testing with TCP ACK-pairs. In *Proceedings of the 7th ACM SIGCOMM conference on Internet measurement*, pages 365–378. ACM, 2007.
- [50] K. Weiss, T. M. Khoshgoftaar, and D. Wang. A survey of transfer learning. *Journal of Big Data*, 2016.
- [51] J. Xu, J. Fan, M. Ammar, and S. B. Moon. On the design and performance of prefix-preserving IP traffic trace anonymization. In *ACM SIGCOMM*, pages 263–266. ACM, 2001.
- [52] L. Xu, K. Harfoush, and I. Rhee. Binary increase congestion control (BIC) for fast long-distance networks. In *INFOCOM*, volume 4, pages 2514–2524. IEEE, 2004.
- [53] P. Yang, J. Shao, W. Luo, L. Xu, J. Deogun, and Y. Lu. TCP congestion avoidance algorithm identification. *IEEE/ACM Transactions On Networking*, 22(4):1311–1324, 2013.
- [54] Y. Yang, L. Wu, G. Yin, L. Li, and H. Zhao. A survey on security and privacy issues in Internet-of-Things. *IEEE Internet of Things Journal*, 4(5):1250–1258, 2017.
- [55] F. Yarochkin and O. Arkin. Xprobe2- A'Fuzzy' Approach to Remote Active Operating System Fingerprinting, 2002.
- [56] M. Zalewski. p0f: Passive OS fingerprinting tool. *Online at http://lcamtuf.coredump.cx/p0f3*, 2017.
- [57] B. Zhang, T. Zou, Y. Wang, and B. Zhang. Remote operation system detection base on machine learning. In *2009 Fourth International Conference on Frontier of Computer Science and Technology*, pages 539–542. IEEE, 2009.