

# Test

XXX

**Abstract**— [1] Accurate application classification is important and useful to improve network performance. However, with the continuous expansion of network scale and the rapid increase of network users, it is very difficult for the existing application classification methods to accurately identify and classify network applications. Currently, most classification methods are suitable for small-scale data sets and cannot achieve high classification accuracy because of the shallow learning structure and the limited learning ability. The emergence of deep learning technology and software-defined networking (SDN) enables the application classification method to process large-scale data. In this paper, by leveraging the SDN architecture, we present a novel hybrid deep neural network-based application classification method, which achieves high classification accuracy without the manual feature selection and extraction. In the proposed application classification framework, by taking the advantage of the logical centralized control and powerful computing capability, the massive network traffic is easily collected and processed by the SDN controller. The processed data is used to train the hybrid deep neural network, which is composed of the stacked autoencoder and softmax regression layer. The deep flow features can be obtained from the stacked autoencoder automatically instead of the manual feature selection and extraction. The softmax regression layer is used as the classifier to realize the application classification. Finally, simulation results demonstrate that our proposed classification method is effective and gets higher classification accuracy than the support vector machine-based classification method

As well known, accurate network application classification is important and useful for many network activities such as network management, quality-of-service (QoS) provisioning, network security, and intrusion detection. 1 Network application classification classifies different network traffic and divides the traffic into different classes. To achieve accurate network applications, a large number of different application classification methods have been proposed in the past few decades, mainly including port-based approach, payload-based approach, and machine learning-based approach. 2 In current years, the most widely used technology for application classification is the machine learning-based approach, in which backpropagation (BP) neural network, Bayesian network, support vector machine (SVM), and C4.5 decision tree are usually applied to the classification. 3-11 Different from identifying port number in the port-based approach as well as inspecting payload content in the payload-based approach, in machine learning-based classification methods, the flow features statistics, including the size of packet, interpacket time, and flow duration time, are used to automatically identify and classify network application by using machine learning. The shallow neural network is generally used to build the application classifier in lots of machine learning-based classification methods. The shallow neural network has limited feature learning ability because of few nonlinear feature extraction layers, and it is more suitable to deal with small-scale data problems. However, with the constant enlargement of network scale and the coming of the big data era, network traffic explosively grows, and large amounts of novel network applications have generated. It is very challenging for traditional shallow neural network-based classification method to deal with massive network traffic because of the limited feature learning ability. In order to extract deep feature and obtain high-level feature presentation from a large-scale data set, deep learning network is required. Recently, deep learning has been proposed, and some deep neural network models, for example, stacked autoencoder, convolutional neural network (CNN), and deep belief network (DBN), have been widely applied in many fields such as classification, speech recognition, and natural language processing. 12 Deep learning is a novel machine learning algorithm that is used to train deeper neural networks. By contrast, traditional learning algorithm, which is adopted to train shallow neural networks, is not suitable for training deep neural networks because of the reason of easily getting stuck in local optima and diffusion of gradient. 13 Compared with traditional shallow neural networks, deep neural networks contain multiple hidden layers such that it has the stronger nonlinear feature extraction ability to achieve the higher level or more complex feature presentation. Meanwhile, deep neural networks have the ability of yielding higher classification accuracy than shallow neural networks. Hence, in this paper, we propose a novel classification method on the basis of the hybrid deep learning network model.

**Index Terms**—Application classification, Encrypted traffic classification, LSTM, CNN



## 1 INTRODUCTION

[2] With a profusion of network applications, traffic classification plays a crucial role in network management and policy-based security control. The widely used encryption transmission protocols, such as the secure socket layer/transport layer security (SSL/TLS) protocols, lead to the failure of traditional payload-based classification methods. Existing methods for encrypted traffic classification cannot achieve high discrimination accuracy for applications with similar fingerprints. In this paper, we propose an attribute-aware encrypted traffic classification method based on the second-order Markov Chains. We start by exploring approaches that can further improve the performance of existing methods in terms of discrimination accu-

racy, and make promising observations that the application attribute bigram, which consists of the certificate packet length and the first application data size in SSL/TLS sessions, contributes to application discrimination. To increase the diversity of application fingerprints, we develop a new method by incorporating the attribute bigrams into the second-order homogeneous Markov chains. Extensive evaluation results show that the proposed method can improve the classification accuracy by 29

NETWORK traffic is composed of packets carrying data belonging to a variety of applications. Classification of traffic helps network operators to identify specific applications and protocols that exist in a network, which can be useful for

many different purposes [19], [32], [33], such as network planning, application prioritization for QoS guarantees, and policy deployment for security control. For instance, a network operator might want to assign traffic from known popular applications with a higher priority for better user experience, and an enterprise network operator might block traffic from a given application by means of application-level firewalls. Traditional traffic classification techniques are often designed based on the analysis of packet contents, such as the port-based methods and the payload-based methods. In recent years we have seen a dramatic growth in the usage of encryption protocols, such as the SSL/TLS protocols [13], [16]. Since packet payloads are encrypted, these methods can no longer fulfil efficient recognition. Therefore, it is desirable to develop classification methods for encrypted traffic [12]. Korczynski and Duda proposed such a method by using the

[3] A Network Traffic Classifier (NTC) is an important part of current network monitoring systems, being its task to infer the network service that is currently used by a communication flow (e.g. HTTP, SIP...). The detection is based on a number of features associated with the communication flow, for example, source and destination ports and bytes transmitted per packet. NTC is important because much information about a current network flow can be learned and anticipated just by knowing its network service (required latency, traffic volume, possible duration...). This is of particular interest for the management and monitoring of Internet of Things (IoT) networks, where NTC will help to segregate traffic and behavior of heterogeneous devices and services. In this paper, we present a new technique for NTC based on a combination of deep learning models that can be used for IoT traffic. We show that a Recurrent Neural Network (RNN) combined with a Convolutional Neural Network (CNN) provides best detection results. The natural domain for a CNN, which is image processing, has been extended to NTC in an easy and natural way. We show that the proposed method provides better detection results than alternative algorithms without requiring any feature engineering, which is usual when applying other models. A complete study is presented on several architectures that integrate a CNN and an RNN, including the impact of the features chosen and the length of the network flows used for training.

[4] Much work has been done on analysing traffic from workstations and web browsers [4]. At first glance, fingerprinting smartphone apps may seem to be a simple translation of existing work. While there are some similarities, there are nuances in the type of traffic sent by smartphones and the way in which it is sent that makes traffic analysis on smartphones distinct from traffic analysis on traditional workstations [5]–[8]. We outline related work by first enumerating traffic analysis approaches on workstations (Section II-A), and then focusing on traffic analysis on smartphones (Section II-B)

[5]

[6] Network traffic classification is defined as a classification of the network flows that are a mixture of various applications with different application protocols.<sup>1</sup>(Gomes JV, Inácio PRM, Pereira M, et al. Detection and classification of peer-to-peer traffic: a survey. ACM Computing Surveys.

2013;45(3):1–40.) It is the foundation of high-performance network protocol design and network operation management.

2.1 — Based on port numbers Traditional methods based on port numbers traffic classifiers simply inspect TCP or UDP port numbers and identify the application layer protocols according to the Internet Assigned Numbers Authority (IANA) list of well-known ports and registered ports.<sup>9,22</sup> The method was simple and fast in the past. However, it has become obsolete nowadays. The mapping between the ports and the target applications is getting more and more blurred. Thereby port numbers as a classification mechanism has not been applied, and it is difficult to deploy.<sup>23,24</sup>

2.2 — Based on deep packet inspection Deep packet inspection methods, usually the most accurate, are based on inspection of the packets' payload. They rely on a database of previously known signatures that are associated to application protocols and search each packet for strings that match any of the signatures.<sup>1,5,6</sup> Searching feature string by DPI is generally in the application layer, which is the load of TCP or UDP. Nevertheless, the main drawbacks of DPI techniques are the following: (1) There are more and more nonstandard applications and private protocols without the open and available protocol specification. This makes the feature string vary and hard to find. (2) Protocol syntax or semantic analysis of data needs strong computation power, leading to a great system overhead.<sup>25</sup> So today, DPI is generally used in traffic identification of the specific applications or as a supplementary means of tagging the network dataset.

2.3 — Based on protocol analysis Based on open protocol regulations, protocol analysis methods analyze the protocols using the following 3 ways for traffic classification: (1) establishing protocol state machines, (2) using fingerprint (protocol traffic features or behavior features) mining, and (3) analyzing flow features and behavior features of unknown protocols using software conversation approach. However, it is very difficult to resolve and obtain effective features because of the nonstandard applications and encryption protocols, degrading the classification quality.

2.4 — Based on machine learning techniques The machine learning methods capture and identify the traffic data packets on the basis of calculating the statistical information of the specific application traffic. The methods use various machine learning algorithms, including supervised and unsupervised learning algorithms. Supervised learning builds a classification model from a training set of labeled instances, which is then used to classify unknown instances. Alternatively, unsupervised learning groups instances that

Network Traffic Classifier (NTC) is an important part of current network monitoring systems, being its task to infer the network service that is currently used by a communication flow (e.g. HTTP, SIP...). The detection is based on a number of features associated with the communication flow, for example, source and destination ports and bytes transmitted per packet. NTC is important because much information about a current network flow can be learned and anticipated just by knowing its network service (required latency, traffic volume, possible duration...). This is of particular interest for the management and monitoring of Internet of Things (IoT) networks, where NTC will help to segregate traffic and behavior of heterogeneous devices

and services. In this paper, we present a new technique for NTC based on a combination of deep learning models that can be used for IoT traffic. We show that a Recurrent Neural Network (RNN) combined with a Convolutional Neural Network (CNN) provides best detection results. The natural domain for a CNN, which is image processing, has been extended to NTC in an easy and natural way. We show that the proposed method provides better detection results than alternative algorithms without requiring any feature engineering, which is usual when applying other models. A complete study is presented on several architectures that integrate a CNN and an RNN, including the impact of the features chosen and the length of the network flows used for training.

In the recent times, the classification of internet traffic data has become a topic with an increased popularity, and its importance has also increased. In particular, the classification of traffic data in the corporate computer networks allows for users to use the network with a good performance and provides opportunities for the advantages such as increasing the quality of service provided to the network administrator, being able to prevent the unwanted traffics, being able to detect the network intrusion attempts and creating the desired security policies on the network

[7] As a fundamental tool for network management and security, traffic classification has attracted increasing attention in recent years. A significant challenge to the robustness of classification performance comes from zero-day applications previously unknown in traffic classification systems. In this paper, we propose a new scheme of Robust statistical Traffic Classification (RTC) by combining supervised and unsupervised machine learning techniques to meet this challenge. The proposed RTC scheme has the capability of identifying the traffic of zero-day applications as well as accurately discriminating predefined application classes. In addition, we develop a new method for automating the RTC scheme parameters optimization process. The empirical study on real-world traffic data confirms the effectiveness of the proposed scheme. When zero-day applications are present, the classification performance of the new scheme is significantly better than four state-of-the-art methods: random forest, correlation-based classification, semi-supervised clustering, and one-class SVM. Traffic classification is fundamental to network management and security [1], which can identify different applications and protocols that exist in a network. For example, most QoS control mechanisms have a traffic classification module in order to properly prioritize different applications across the limited bandwidth. To implement appropriate security policies, it is essential for any network manager to obtain a proper understanding of applications and protocols in the network traffic. Over the last decade, traffic classification has been given a lot of attention from both industry and academia. There are three categories of traffic classification methods: port-based, payload-based, and flow statistics-based [2]. The traditional port-based method relies on checking standard ports used by well-known applications. However, it is not always reliable because not all current applications use standard ports. Some applications even obfuscate themselves by using the well-defined ports of other applications. The payload-based method searches for

the application's signature in the payload of IP packets that can help avoid the problem of dynamic ports. Hence, it is most prevalent in current industry products.

[8] Network communication became the standard way of exchanging information between applications located on different hosts. The exchanged application-layer data is segmented and encapsulated into IP packets, which are transmitted through the network. Deep Packet Inspection (DPI) tools analyze the content of the packets by searching for specific patterns (i.e., signatures). Thus, DPI became one of the fundamental traffic analysis methods for many tools performing traffic classification, network management, intrusion detection, and network forensics.

The traffic classification process labels traffic based on a predefined set of classes (e.g., applications). A commonly used classification approach is Deep Packet Inspection (DPI), which classifies traffic by inspecting the payload of each packet traversing the inspection point. DPI attains remarkably high accuracy for unencrypted traffic. Nonetheless, since DPI relies on the visibility of the payload to the classifier, its effectiveness diminishes with usage of data encryption, which is the conventional wisdom, or when examining the content of the packets traversing the network is forbidden or limited due to privacy or complexity concerns.

[7]. For example, most QoS control mechanisms have a traffic classification module in order to properly prioritize different applications across the limited bandwidth. To implement appropriate security policies, it is essential for any network manager to obtain a proper understanding of applications and protocols in the network traffic. Over the last decade, traffic classification has been given a lot of attention from both industry and academia. There are three categories of traffic classification methods: port-based, payload-based, and flow statistics-based [2]. The traditional port-based method relies on checking standard ports used by well-known applications. However, it is not always reliable because not all current applications use standard ports. Some applications even obfuscate themselves by using the well-defined ports of other applications. The payload-based method searches for the application's signature in the payload of IP packets that can help avoid the problem of dynamic ports. Hence, it is most prevalent in current industry products. However, more often than not, the payload-based method fails with encrypted traffic. In recent academic research, significant attention has been given to applying machine learning techniques to the flow statistics-based method. The statistical method only uses flow statistical features, such as interpacket time, without requiring deep packet inspection (DPI).

A flow consists of successive IP packets with the same 5-tuple: source IP, source port, destination IP, destination port, transport protocol.

## 2 RELATED WORK

[1] In recent years, many research works have already applied machine learning methods in network application classification. [2-11, 23, 24] Most of them are focused on improving the machine learning algorithm and feature selection since the selected features and machine learning algorithm have a great effect on the classifier's performance.

However, with the explosive growth of network traffic, it is very difficult for the current network architecture to deal with large amounts of data. Therefore, to obtain a set of optimal flow features and improve the classification accuracy, Santos da Silva et al [24] designed a novel flow features selection architecture, which could determine the optimal subset of flow features for the classification of different types of traffic flows by taking advantage of the SDN architecture. Different from most research works on application classification, to achieve the classification and satisfy the QoS requirement of the network application at the same time, Wang et al [11] devised a QoS-aware traffic classification framework in SDN. In this framework, deep packet inspection technology was used to detect the elephant flow, and the semisupervised machine learning algorithm was used to achieve the QoS-aware traffic classification through the mapping function. Specifically, the application flow was mapped to a certain predefined QoS class according to flow features in the mapping function. Inspired by the aforementioned research works, in this paper, we propose an application classification framework by combining SDN and deep learning. With the powerful computing capability, we use the controller to deal with the massive network traffic and flow statistics. We construct a hybrid deep learning network model, which is composed of the stacked autoencoder and the softmax regression layer, to extract flow features and build the application classifier [5].

The application of machine learning for the detection of malicious network traffic has been well researched over the past several years because traditional pattern-matching approaches cannot be used. Unfortunately, the promise of machine learning has been slow to materialize in the network security domain. In this paper, we highlight two primary reasons why this is the case: inaccurate ground truth and a highly non-stationary data distribution. To demonstrate and understand the effect that these pitfalls have on experiments that show how six common algorithms perform when confronted with real network data.

[9] Botnets represent one of the most aggressive threats against cyber security. Different techniques using different feature sets have been proposed for botnet traffic analysis and classification. However, no work has been performed to study the effect of such differences. In this paper, we perform a study on the effect of (if any) the feature sets of network traffic flow exporters. To this end, we explore five different traffic flow exporters (each with a different set of flow features) using two different protocol filters [Hypertext Transfer Protocol (HTTP) and Domain Name System (DNS)] and five different classifiers. We evaluate all these on eight different botnet traffic data sets. Our results indicate that the use of a flow exporter and a protocol filter indeed has an effect on the performance of botnet traffic classification. Experimental results show that the best performance is achieved using Tranalyzer flow exporter and HTTP filter with the C4.5 classifier.

[10]

Network traffic classification is a critical network processing task for network management. Traffic measurement and classification enable network administrators to understand the

current network state and reconfigure the network such that the observed network state can be improved over time. The complexity and dynamic characteristic of today's network traffic have necessitated the need for traffic classification techniques that are able to adapt to new concepts. This includes the ability to classify types of traffic almost instantaneously to avoid outdated the knowledge gained from the learning of new concepts. Data stream mining techniques are able to classify evolving data streams such as network traffic in the presence of concept drift. In order to classify high bandwidth network traffic in real-time, data stream mining classifiers need to be implemented on reconfigurable high throughput platform, such as Field Programmable Gate Array (FPGA). This paper proposes an algorithm for

-means classification using Manhattan distance can classify network traffic 3 times faster than Euclidean distance at 671 thousands flow instances per second

[11] A Survey of Payload-Based Traffic Classification Approaches

[12] A survey of methods for encrypted traffic classification and analysis

### 3 METHODOLOGY

First, we would like to introduce the definition of a network flow in a TCP/IP-based network as mainly used in this work. It is a five-dimensional tuple having the following fields:  $\{$ source IP, destination IP, source port, destination port, transport protocol  $\}$ . The transport protocol could be either TCP or UDP. A flow is basically a client-server communication or a dialog between two peers. We consider this concept important as it is close to the traffic patterns in real networks that we want to classify more accurately. For example, an HTTP client-server communication is a bidirectional flow that is initiated by the client. Authors in [7] analyze the traffic patterns inside a data center and define the concept of flow. They find that more than 80% of center flows last less than 10 s.

### 4 EXPERIMENTS

In this paper, we used the Andrew Moore 28 datasets, which consisted of 10 separate subdatasets each from a different period of the 24-hour day. The day trace was split into 10 blocks of approximately 1680 seconds. Each subdataset was represented by a data text file that included tens of thousands of data lines. Each line represented a traffic flow. The information is derived from packet header information. To reduce the imbalance of the data, we deleted the traffic flows of games and interactive, which the samples were very few. And we extracted samples randomly within 3000 from every subset to build the new datasets. The datasets included 24 897 samples as in Table 1. The generation process of the training set in subsequent experiments is as follows:

To evaluate the performance of the proposed application classification method, we have conducted extensive simulation experiments.

**TABLE 2** Network application classification<sup>8</sup>

Classification	Application
Bulk	ftp
Database	postgres, sqlnet oracle, ingress
Interactive	Ssh, klogin, rlogin, telnet
Mail	imap, pop2/3, smtp
Services	X11, dns, ident, ldap, ntp
WWW	www
P2P	KaZaA, BitTorrent, GnuTella
Attack	Internet worm and virus attacks
Games	Microsoft Direct Play
Multimedia	Windows Media Player, Real

Abbreviations: P2P, peer-to-peer; WWW, World Wide Web.

Fig. 1. Simulation results for the network.

#### 4.1 Dataset

[] We evaluate the proposed deep learning-based application classification method on the real-world traffic data set. We select the Moore data set for testing the application classification method, which is obtained from the computer laboratory in the University of Cambridge and has been widely used in many traffic classification research works. 2,6,8,23,34 The data set is composed of 10 separate data sets collected in the different period of a day, and each data set is only composed of TCP traffic flows. In each data set, for any TCP traffic flow, its 248 flow features such as the size of flow and the duration time of flow and the corresponding application class label are recorded. All network flows are categorized into 10 classes (ie, World Wide Web (WWW), Mail, Bulk, Services, peer-to-peer (P2P), Database, Attack, Interactive, Games, and Multimedia). The traffic classification and the corresponding network applications are summarized in Table 2. Furthermore, to keep the balance of the sample data set, we delete the corresponding 2 types of application samples from the 10 data sets because the sample numbers of Games and Interactive are relatively small. Moreover, to make the data set more uniform and to get more accurate simulation results, we randomly select sample data from the 10 data sets and classify all the network applications into 10 classes (ie, WWW, Mail, File Transfer Protocol (FTP)-control, FTP-pasv, Attack, P2P, Database, FTP-data, Multimedia, and Services).

## 5 DISCUSSIONS

### 5.1 Compare with related work

### 5.2 Limitations

## 6 CONCLUSION

The conclusion goes here 6.

[13] [4] [2]

## APPENDIX A

### PROOF OF THE FIRST ZONKLAR EQUATION

Appendix one text goes here.

## ACKNOWLEDGMENTS

The authors would like to thank...

## REFERENCES

- [1] C. Zhang, X. Wang, F. Li, Q. He, and M. Huang, "Deep learning-based network application classification for SDN," *Transactions on Emerging Telecommunications Technologies*, vol. 29, no. 5, 2018.
- [2] M. Shen, M. Wei, L. Zhu, and M. Wang, "Classification of Encrypted Traffic with Second-Order Markov Chains and Application Attribute Bigrams," *IEEE Transactions on Information Forensics and Security*, vol. 12, no. 8, pp. 1830–1843, 2017.
- [3] M. Lopez-Martin, B. Carro, A. Sanchez-Esguevillas, and J. Lloret, "Network Traffic Classifier with Convolutional and Recurrent Neural Networks for Internet of Things," *IEEE Access*, vol. 5, pp. 18 042–18 050, 2017.
- [4] V. F. Taylor, R. Spolaor, M. Conti, and I. Martinovic, "Robust Smartphone App Identification via Encrypted Network Traffic Analysis," *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 1, pp. 63–78, 2018.
- [5] B. Anderson and D. McGrew, "Machine Learning for Encrypted Malware Traffic Classification," *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining - KDD '17*, vol. Part F1296, pp. 1723–1732, 2017. [Online]. Available: <http://dl.acm.org/citation.cfm?doid=3097983.3098163>
- [6] J. Cao, Z. Fang, G. Qu, H. Sun, and D. Zhang, "An accurate traffic classification model based on support vector machines," *International Journal of Network Management*, vol. 27, no. 1, pp. 1–15, 2017.
- [7] J. J. Zhang, X. Chen, Y. Y. Xiang, W. W. Zhou, and J. J. Wu, "Robust Network Traffic Classification," *TNetworking*, vol. 23, no. 4, pp. 1–14, 2014.
- [8] T. Bujlow, V. Carela-Español, and P. Barlet-Ros, "Independent comparison of popular DPI tools for traffic classification," *Computer Networks*, vol. 76, pp. 75–89, 2015. [Online]. Available: <http://dx.doi.org/10.1016/j.comnet.2014.11.001>
- [9] F. Haddadi and A. N. Zincir-Heywood, "Benchmarking the Effect of Flow Exporters and Protocol Filters on Botnet Traffic Classification," *IEEE Systems Journal*, vol. 10, no. 4, pp. 1390–1401, 2016.
- [10] H. R. Loo, S. B. Joseph, and M. N. Marsono, "Online Incremental Learning for High Bandwidth Network Traffic Classification," *Applied Computational Intelligence and Soft Computing*, vol. 2016, 2016.
- [11] M. Finsterbusch, C. Richter, E. Rocha, J. A. Müller, and K. Hänßen, "A survey of payload-based traffic classification approaches," *IEEE Communications Surveys and Tutorials*, vol. 16, no. 2, pp. 1135–1156, 2014.
- [12] I. G. Siqueira, L. B. Ruiz, and A. Loureiro, "A survey of methods for encrypted traffic classification and analysis," *International Journal of Network Management*, no. October 2005, pp. 17–31, 2014. [Online]. Available: <http://onlinelibrary.wiley.com/doi/10.1002/nem.604/abstract>
- [13] W. Wang, Y. Sheng, J. Wang, X. Zeng, X. Ye, Y. Huang, and M. Zhu, "HAST-IDS: Learning Hierarchical Spatial-Temporal Features Using Deep Neural Networks to Improve Intrusion Detection," *IEEE Access*, vol. 6, pp. 1792–1806, 2017.