

Deep RNN-based Traffic Analyzation Scheme for Detecting Target Applications

Author 1, Author 2, Author 3, Author 4 and Author 5

Abstract—This letter proposes a deep RNN-based traffic alayzation scheme for detecting target applications. The proposed scheme improves the classification performance in telecommunication networks with heavy traffic. In this work, we design a novel classification learning method where input features and output labels are two-dimensional image with traffic packet and target applications, respectively. Specifically, traffic packets are provided through universitat politecnica de catalunya barcelonatech as the training data[1]. Then, the proposed scheme provides a fast and exact traffic alayzation from the produced inferred function by analyzing the training data. We implement the proposed scheme using a commercial deep long short-term memory system. Simulation-based experiments show that the proposed scheme achieves almost xx% accuracy performance with low complexity, requiring only xx ms elapsed time.

Index Terms—LSTM, Deep Learning, Network Traffic, Traffic Analyzation.

I. INTRODUCTION

VARIOUS types of applications and services are operating in recent networks. It also provides easy access to the network from anywhere, including the smartphone, laptop, and desktop users. As the quality of demand increases, larger and more diverse traffic data is occurring. Therefore, the scope of human data analysis and management is greatly expanding[2]. Therefore, it is necessary to establish a lightweight and automated network[3]. In addition, the application of Deep Learning, which is gaining attention recently, is increasing. Some of deep learning technologies include Multilayer Perceptron (MLP), CNN, and Recurrent Neural Network (RNN). The MLP is a basic deep learning model consisting of approximately one or more hidden layers. CNN, which is a neural network that is comprised of three levels of local reception area, convolutional layer, and pooling layer, is a deep learning technique that is typically used in image analysis. RNN is a circular neural network suitable for processing random sequential data, and its characteristics are influenced by previous computational results. In this paper, RNN is used among deep learning techniques to classify network traffic based on flow. The flow is a collection of the same traffic by 5-tuple, which in turn contains network packets. Therefore, RNN, which is suitable for training sequential data, is used for flow-based network traffic classification. This requires a process of converting flow data collected on the network into a form that can be learned using RNN. To this end, through the traffic data preprocessing process, the learning data generation step makes the raw traffic data divided by each application into traffic split and splited network traffic set easy to learn. The RNN learning will then categorize the flow. Chapter 2 of this paper describes the system model, and Chapter 3 suggests deep RNN-based traffic analyzation scheme, and Chapter 4 conducts

performance evaluation. Chapter 5 describes conclusions and future research plans.

II. SYSTEM MODEL

We need to express a system model that is behind the proposed scheme.

III. PROPOSED DEEP RNN-BASED TRAFFIC ANALYZATION SCHEME (DR-TAS)

This chapter describes traffic data preprocessing and deep RNN model for Proposed deep RNN-based traffic analyzation scheme.

A. traffic data preprocessing

The traffic data used to perform traffic classifications is a preprocessing packet capture (PCAP) files supplied by universitat politecnica de catalunya barcelonatech (UPC) [1] suitable for RNN learning. PCAP is a file that captures network packets and stores them in the form of PCAP files using programs such as Wireshark and TCPdump. PCAP provided for flow-based network traffic classification using RNN suggested in this letter need to be pre-processed into data sets classifiable in RNN traffic analyzation. This requires a data pre-processing process to filter and shorten the PCAP files. The original PCAP file is approximately 59GB in size and has 769507 flows. In the packet_default.info file provided with the PCAP file, there is labeling for the traffic data. Labeling is divided into classes, such as the application type, protocol, and application name for the corresponding traffic flow. The labeling file gave us the correct label for ground truth in the flow-based traffic classification with RNN suggested in this letter.

For the data pre-processing process, eight types of applications were selected based on the number of flows. The top eight label names were Remote Desktop Protocol (RDP), Skype, SSH, Bittorrent, HTTP-face-Govert, HTTP-Doki, HTutube, and only the corresponding label selection flow. Only the payload of the application layer was filtered out of the selected flow internal packet, and the extracted payload data was placed and eight application layer Payload data files were generated. And for HTTP-facebook-Google, HTTP-Web, HTTP-Doki, HTTP-Youtube, we finally combined the label with the label Web to create a total of five data sets. Through the integration process, the application layer Payload data files of RDP, Skype, SSH, Bittorrent, and Web were completed. The data file has header information for each flow, and the payload separated by a # delimiter for each packet. That is, each file is a collection of flows with the same label, with packets having only a payload and the payload being separated by a packet. The payload is also a string, with a size of four bits per one

TABLE I
STATISTICAL INFORMATION OF EACH APPLICATION

	rdp	skype	ssh	bittorrent	http-web
Total number of filtered flows	153,349	2,041	38,831	96,222	21,715
Total number of packets(K)	6,876K	3,015K	173,911K	207,306K	47,136K
Average number of packets in a flow	44	1,477	461	2,154	2,483
Total packet size (GB)	131.3	12.0	321.7	2,170	698.4
Average packet size over all flows (KB)	20.0	4.2	18.8	11.3	15.9
Min packet size over all flows (B)	8.0	8.0	8.0	8.0	8.0
Max packet size over all flows (MB)	5.5	0.3	6.6	0.8	40.0

character. Figure 1 shows statistical information on the finished application layer Payload data file.

This step describes the process of converting LSTM into learning data using the application layer payload data generated in the previous step. A data set for learning will have flight data with 8,750 flows, 1,250 validation data, and 2,500 test data. Each label data equal to the number of each transaction, validation, test data set. Each application flow in the application layer payload data file goes through a conversion process into data for LSTM learning. The conversion job removes the head information that indicates the flow id, start time, and end time of each flow, and imports the data portion of the application layer only the payload in the flow unit and creates a data set. The LSTM models should have the same number of data input at a time to perform the learning in flow units. That is, a flow should have the same size packets. For this purpose, the user can specify as many packets per flow as P. In addition, one packet in the flow has a user-designed size, which can be set by the user to N, and it has a one-dimensional array size. The pixels of a packet's pixel have a 2^n bit size and a character-type value replaced by a floating point value. That is, the data in one flow will be the same size as the equation (1).

$$N_{\text{flow}} = N \times 2^n \times P \quad (1)$$

Fig.1 shows the payload of a packet in one of the completed flows in an two-dimentional image, a element of four bits in size, and a value between zero and fifteen floating point. In addition, due to the nature of the LSTM network structure, all flows should have the same number of packets.Because the length of sequence is set in advance, the number of packets per flow should be the same by that length. However, since not all Applicatoin flows can have the same number of packets, if the number of packets per flow is smaller than the N set, then the packet is generated by zero. Train, Validation and Test Label have labels for a total of five applications, which can be expressed as one-hot vectors of five lengths. A one-hot vector label having a value of only one element, can be defined as a label on which an index of one points to an application name. Table 2 shows the structure represented by the Train, Validation, and Test labels as on-hot vectors.

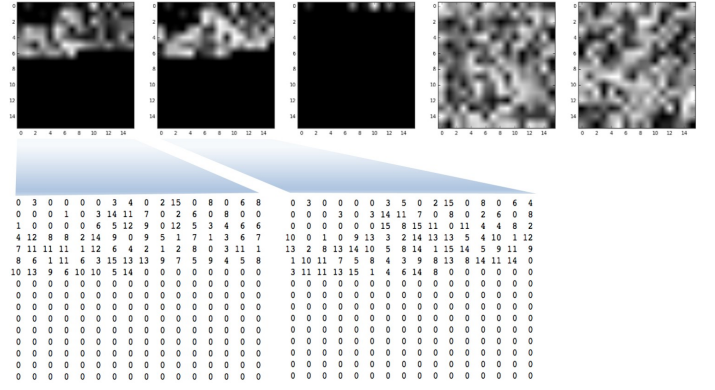


Fig. 1. Input features of flow packet

TABLE II
OUTPUT LABELS OF APPLICATION ONE-HOT VECTOR

Application	RDP	Skype	SSH	Bittorrent	Web
label	[1,0,0,0,0]	[0,1,0,0,0]	[0,0,1,0,0]	[0,0,0,1,0]	[0,0,0,0,1]

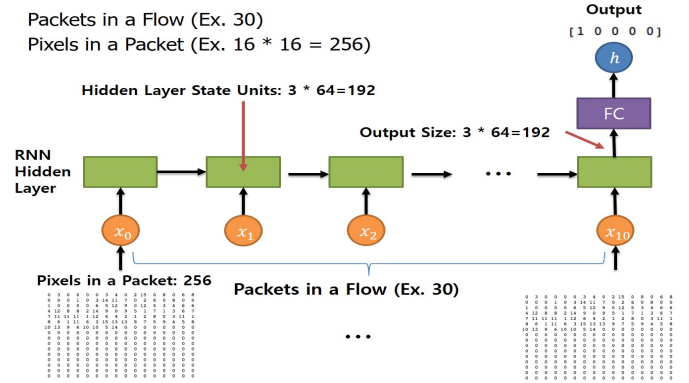


Fig. 2. Proposed deep RNN-based traffic analysis scheme

B. Proposed deep RNN model

This subsection describes the LSTM models that will be used for flow based network traffic classification using the proposed LSTM. LSTM shows great performance for various natural language processing (NLP) issues, especially LSTM is able to learn sequential data over CNN, which is used more by Deep Learning. Therefore, we believe that the network traffic classification through LSTM will be accurate in analyzing the flow that contains sequential information of packets. Fig.2 shows the learning structure of the Proposed deep RNN-based traffic analysis scheme used to classify traffic, and the learning model consists of a single layer. The time step for learning LSTM is to preset the number of packets throughout the flow to be learned. Time steps represent the number of inputs a flow enters sequentially from a single layer to an LSTM. Therefore, the form of learning data sets is represented in the train data set (number of flows, number of packets per flow, and payload size per packet), and the label data (number of flows, number of labels). The data sets of tests and validation are also inputted in the same form and learned through LSTM. The problem with the existing vanilla RNN is that of long-term dependency if the length of sequence is longer. Thus, this network architecture utilizes LSTM, a solution to long-term dependency problems.

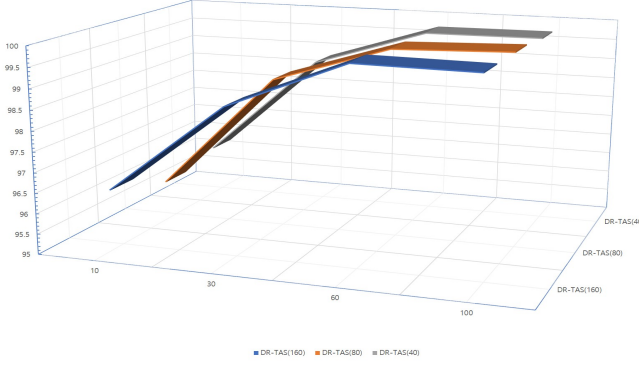


Fig. 3. Accuracy of deep RNN-based traffic analysis scheme

IV. PERFORMANCE EVALUATION

In this section, we present the results of classification of traffic data through experiments of LSTM model that classify target application. We also compare the network traffic classification results of LSTM model and CNN model. The experiment is divided into 4 kinds, and the first is the experiment which shows the accuracy of LSTM. The second is an experiment to compare the accuracy of each application of LSTM. To compare the accuracy of each application, we use precision, recall, and f1-score as a measure of comparison. The third is an experiment to compare CNN and LSTM accuracy. We compare the accuracy of CNN and LSTM by increasing the number of applications to 2, 3, 4, and 5 for comparison. Finally, it is an experiment to compare the elapsed time in the third experiment.

A. Experiment Environment

Experiments were conducted on Ubuntu 16.04 LTS, using 32GB of RAM and two NVIDIA GTX 1080Ti 11GB. We also used Tensorflow developed by Google Brain Team in Python 3.6 environment to configure LSTM and CNN deep learning model. For the experiment, we have 2000 flows for each application in train data, and the number of packets per flow is set to 10, 30, 60, and 100. In addition, the payload size of each packet was set to 40, 80, and 160 to generate a total of 12 train data.

B. LSTM Model Evaluation

In the LSTM model evaluation experiment, it is an experiment to examine the performance of the LSTM model. This is an experiment that compares accuracy based on the number of packets per flow and the size of each payload. In figure 3, the x-axis is the number of packets per flow and the y-axis is the accuracy. As a result, it can be confirmed that the more the number of packets per flow, the higher the accuracy. In addition, it shows that the accuracy is high as the payload size of each packet increases.

C. LSTM Model Evaluation for Each application

The performance of LSTM model for each application is compared. Confusion Matrix is used to compare performance of each application of LSTM model. The Confusion Matrix is

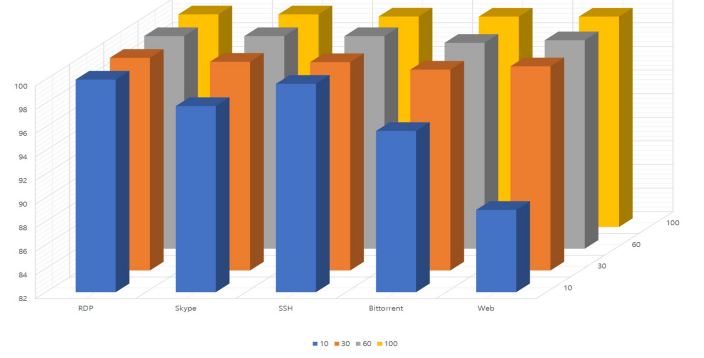


Fig. 4. Accuracy of DR-TAS for each target application

the most widely used method for representing binary classification evaluation results. The Confusion Matrix shows you how well your deep-running model is categorized for each application. Figure x shows the Confusion Matrix, where the rows represent the negative class and the positive class, and the columns represent the predicted by negative and the predicted by positive. For example, TN (True Negative) refers to the case where the prediction result of the model is a negative class, which is actually a negative class. FP (False Positive) means that the model is predicted as a positive class, but actually it is a negative class. FN (False Negative) indicates that the model is predicted as a negative class, but the actual result indicates a positive class, and TP (True Positive) indicates a case where both the predicted result and the actual result are positive. Calculate precision, recall and f1-score based on Confusion Matrix.

The precision means the percentage of actual positive predictions out of the positive predicted class, and the equation of precision is (2).

$$precision = \frac{TP}{TP + FP} \quad (2)$$

The recall is called the positive detection rate, and it means the percentage of the total positive class predicted by a positive class, and the equation of recall is (3).

$$recall = \frac{TP}{TP + FN} \quad (3)$$

Finally, f1-score means the harmonic mean of precision and recall, and the equation is (4).

$$f1 - score = \frac{2 \times precision \times recall}{precision + recall} \quad (4)$$

The recall and precision are in conflict, and the f1-score is an index that shows the accuracy at once by integrating precision and recall.

Figure 4 shows the accuracy of each application, and figure 5 shows the f1-score value for each application. In the case of a graph showing accuracy, the accuracy of each application increases as the number of packets per flow increases. Also in the case of f1-score, it can be seen that the value increases as the number of packets per flow. However, in the case of web and bittorrent, the smaller the number of data, the lower the accuracy and the f1-score value than other applications.

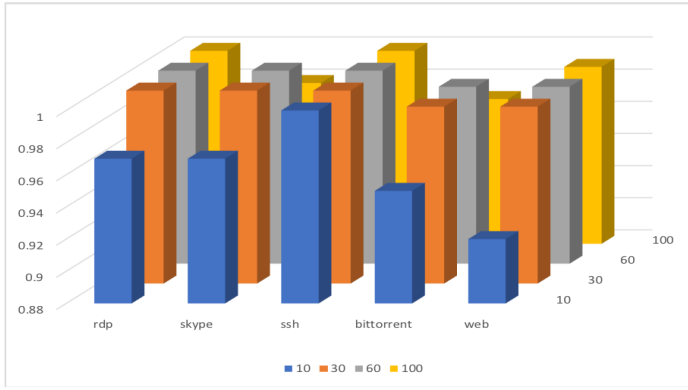


Fig. 5. f1-score of DR-TAS for each target application

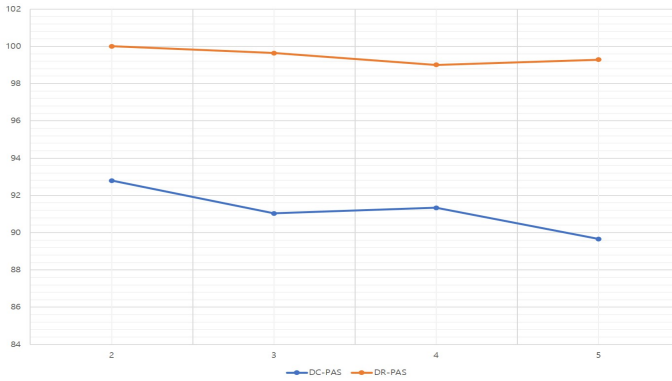


Fig. 6. Comparison on accuracy of DR-TAS vs. DC-TAS

D. Comparison of CNN and LSTM performance

This is a section to compare the performance of CNN and LSTM deep learning model. we compare the accuracy of each application by increasing it to 2, 3, 4, or 5 classes. Two applications compared CNN and LSTM by selecting RDP and Skype with high f1-score. In three cases, ssh was added to the two applications selected above. Bittorrent was added in 4 cases, and finally web was added to compare the total of 5 accuracies. Figure 5 shows the accuracy comparison results for the number of applications of CNN and LSTM. Both CNN and LSTM show that accuracy is lowered when Bittorrent and web are added. However, the overall accuracy of the LSTM is higher than that of CNN.

E. Comparison of CNN and LSTM elapsed time

We compare elapsed time from data preprocessing of CNN and LSTM to result prediction. The elapsed time is the sum of the preprocessing time for making the appropriate data for each model and the time taken to input the training data into the CNN and LSTM model. Figure 6 shows the elapsed time of CNN and LSTM, and shows the time taken to increase the number of applications as above. As the number of applications increases, the elapsed time increases. The elapsed time has increased because the amount of data to learn has increased. However, it can be seen that the elapsed time of 4 and 5 applications that add bittorrent and web data increases dramatically. The reason is that the two application data are larger than RDP, Skype, and ssh, so that the preprocessing process, which makes the data suitable for the model, takes

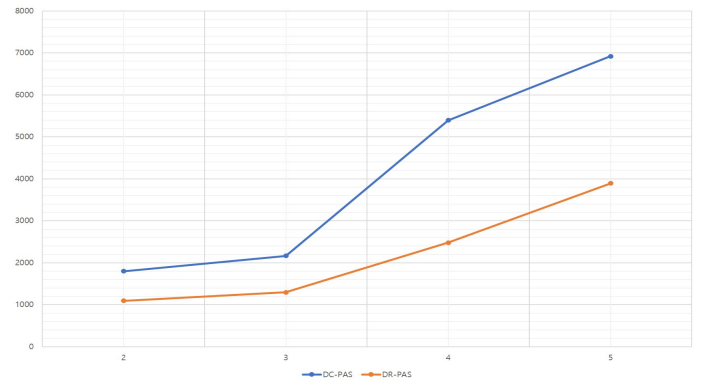


Fig. 7. Comparison on elapsed time of DR-TAS vs. DC-TAS

more time. In addition, bittorrent and web are bigger than other applications because most of the data is image or video. However, in the case of LSTM, the elapsed time is less than half that of CNN as a whole. In the case of CNN, one packet is arbitrarily fetched from each application-specific flow in the preprocessing process, so that it takes a long time to read one flow. On the other hand, in the case of LSTM, since the flow is generated for each application and generated as data, the time of the preprocessing process is shorter than that of CNN.

V. CONCLUSION

As a result of learning from the LSTM model using network traffic data as a unit of flow, its accuracy was over 99%. This confirmed that the classification was nearly 100%. Future research will explore the classification of network traffic through deployment in a real network and also how new packets, other than those already learned, will be classified as they enter the real network.

REFERENCES

- [1] Valentín Carela-Español, Tomasz Bujlow, and Pere Barlet-Ros: "Is Our Ground-Truth for Traffic Classification Reliable?", *In Proc. of the Passive and Active Measurements Conference (PAM'14)*, Los Angeles, CA, USA, March 2014.
- [2] Jinwan Park, "statistics signiture based application traffic classification," *Korea Communication Journal*, vol. 34, pp. 1234-1244, Nov. 2009.
- [3] F. Risso, "Lightweight, Payload-Based Traffic Classification: An Experimental Evaluation," *IEEE International Conference on Communications 2008*, 2008.
- [4] Yann LeCun, "Deep Learning," *Nature International Weekly Journal of Science*,
- [5] T. Mikolov, S. Kombrink, L. Burget, J. Cernocky and S. Khudanpur, "Extensions of recurrent neural network language model," *2011 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5528-5531, 2011.
- [6] H. Sak, Andrew Senior and F. Beaufays, "Long short-term memory recurrent neural network architectures for large scale acoustic modeling," *Proceedings of the Annual Conference of the International Speech Communication Association (INTERSPEECH)*, pp. 338-342, Jan. 2014.