

Network Traffic Classification Model Using Ensemble Machine Learning

Deepali Shukla¹, Dr. Kavi Bhushan², Er. Gur Sharan Kant³

¹Scholar M.Tech, Department of Computer Science & Engineering, Sir Chhotu Ram Institute of Engineering and Technology, Chaudhary Charan Singh University, Meerut, U.P., India

²Assistant Professor Sir Chhotu Ram Institute of Engineering and Technology, Chaudhary Charan Singh University, Meerut, U.P., India

³Assistant Professor Sir Chhotu Ram Institute of Engineering and Technology, Chaudhary Charan Singh University, Meerut, U.P., India

ARTICLE INFO	ABSTRACT
Received: 29 Dec 2024 Revised: 12 Feb 2025 Accepted: 27 Feb 2025	<p>Network Traffic Classification (NTC) is concerned with the identification of different types of application traffic by the examination of the data packets which are received in a communication network. This process is crucial for network management and it is getting more and more important in the last few years. The typical workflow for traffic classification includes data input, data preprocessing, attribute extraction, classification and performance analysis. It is becoming more and more rewarding the role of machine learning techniques in classifying network traffic as they become more and more advanced. This paper presents a novel paradigm intended to increase the efficacy of traffic classification. The suggested model uses a hybrid classification approach along with SMOTE to address class imbalance and Principal Component Analysis (PCA) for feature reduction. This work leverages Python, specifically when it operates in the Anaconda environment, to evaluate the efficacy of the proposed framework.</p> <p>Keywords: Classification, Network traffic, Features, Machine learning</p>

1. Introduction

The fast-paced development of the internet is causing alarm for the protection of people's private lives and networking (security). The popular privacy-preserving tools have been thoroughly tested and services such as Virtual Private Networks (VPNs) and Anonymizing Mechanisms (AMs) developed to protect the confidentiality are widely installed. Proxies, which are the essential constituents of a browser, usually provide users with the means to obscure the information being transferred, thereby they can share the device in a secure way which also helps hide the unneeded data. Consequently, traffic classification stands out as the crucial goal in order to implement QoS, perform network traffic management, and ultimately secure networks [1]. In simple terms, traffic classification is the assignment of the category of traffic moving through a network which is vital for the recognition of diverse threats. Many systems and applications are developed to improve network potential and control, with machine learning (ML) methods being widely employed for this purpose. These techniques aid in the effective distribution, management, and control of network resources. Additionally, a main element of IDS, which are designed to identify risks and malicious activity in networks, is traffic classification.

Nowadays, academia shows an increasing interest in the issue of traffic network identification because of the considerable importance of it [2]. To accomplish this, several ways have been developed and put forward by researchers in the last few decades. The following is a discussion of several methods for classifying network traffic:

a. Port-based classification: This method is implemented to check the TCP or UDP port registered with the IANA for a match with the packet headers and it is supposed to be an application identifier. Some of the newest apps are using unregistered or random port numbers that have no control over them, including peer-to-peer (P2P) applications, which results in an increased chance of false negatives. Moreover, other apps often make their desire to hide their communication clear to avoid detection and surpass filtering that way. They do so because they simply want to dodge the restrictions that the operating system access control mechanisms impose.

b. Payload-based classification: Numerous commercial products and solutions are put out to address the drawbacks and interdependencies of previous methods that prioritize packet headers over contents [3]. This approach is referred to as payload-based classification. This technique examines the packet's contents and compares them to a predetermined signature bank. It often yields more precise and superior outcomes. Furthermore, it is regarded as a crucial initial step in an IDS to detect malicious activities on the network.

c. Statistical classification: This solution is based on a logic-driven method that uses statistical features of flowing traffic to find applications [4]. During the classification process, several flow-level parameters are assessed, including packet presence, packet dimensions, IAT, and inactivity in flowing traffic. These metrics help the classification algorithm differentiate between many applications, as they are distinct and useful for particular application types.

d. Behavioral classification: This classification method entails an examination of all the network traffic sent out by the endpoint. One of the ways to specify an application is by investigating the transmission patterns issued by the target device. Specifically, the TLP and the total number of ports used are the factors that will determine the number of hosts connected. This technique produces top outcomes even on the lowest possible computational cost.

According to Figure 1, the structure of NTC model is presented. Collection of data, attribute identification, attribute reduction and selection, and model creation are the stages that set up this architecture [5]. This process illustrates the stages of NTC approaches that apply machine learning (ML) algorithms to identify and categorize incomprehensible network data.

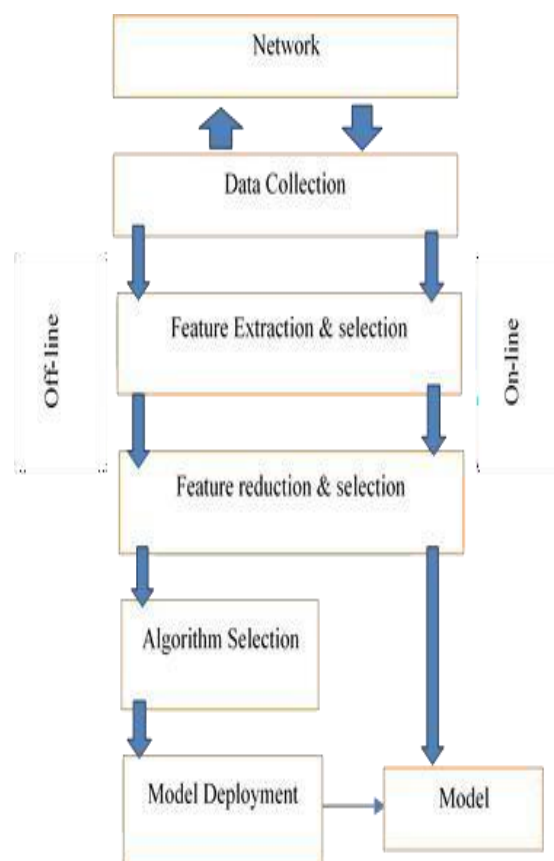


Figure 1: Network Traffic Classification Model

The sections below are an explanation that embodies the stages of the traffic classification model network.

a. Data Collection: Historical data is the basis of machine learning (ML) algorithms. These algorithms' potential and generality are enhanced through the use of exclusively and an extensive dataset that has details of the problem at hand [6]. Categorizations of network traffic processing are handled effectively through this mechanism of classifying

network traffic because of many difficulties, like the interconnectedness and flexibility of web networks, the pace of traffic growth and privacy laws which impose restrictions against data collection. In the data gathering phase, there are many different situations that can be assessed along the network. The purpose of this phase is to gather IP addresses in a certain time frame. Furthermore, it includes many tasks like storage, flow reconstruction, and packet management. While the online method focuses on real-time packet flow, offline processing requires the collection of previous datasets [7].

b. Feature extraction: This phase's primary goal is to extract qualities and document the information that will be used to address the current problem. This stage is crucial for quantifying characteristics, such as data about the state of the process. Additionally, numerous parameters that indicate important attributes in the gathered data are calculated using the attribute extraction technique. The main goal is to obtain descriptors, with the output of this process being a structured table. The table is produced using feature columns and each row represents the current position of every sample, acting as a pattern with an additional random column [8]. When the status is uncertain, the patterns are not given a label.

c. Feature selection and reduction: In order to acquire a smaller collection of new characteristics or to reduce the amount of space needed, attribute selection and reduction procedures are used in this step. This procedure reduces and selects the extracted properties. The selection step seeks to identify a more manageable collection of features that characterize the process, while the reduction process generates new attributes based on the original features. This phase addresses a number of issues, including cost and time consumption. For this goal, three methods are practiced: filter, wrapper, and embedded techniques,

d. Classification: Firstly, the original data is processed into a new data set tailored to the chosen features. After that, the offline process leverages the obtained new dataset to come up with models that do the classification of the specified data. The algorithm selection process involves methods of machine learning model (ML) selection. This method can be seen as a possibility of going through a large number of combinations of algorithms. The identification of the appropriate model is a key point in various machine learning approaches that correctly categorize network data [10].

1.1 Contribution

This work introduces a novel approach for network traffic classification by integrating the SMOTE technique, PCA for dimensionality reduction, and a Voting ensemble classifier. The use of multiple classifiers such as Random Forest (RF), Support Vector Machine (SVM), and K-Nearest Neighbors (KNN) in a voting mechanism improves classification accuracy and reduces the false positive and false negative rates. Additionally, the application of SMOTE addresses class imbalance, ensuring better performance on real-world datasets. The integration of PCA enhances feature selection, reducing computational complexity while maintaining classification effectiveness.

1.2 Objectives

Following are the main objectives of this work:

1. To design and implement an efficient network traffic classification system using the NSL-KDD dataset, focusing on distinguishing between normal and attack traffic.
2. To address class imbalance in the dataset by utilizing the SMOTE technique for generating synthetic minority class samples.
3. To enhance classification accuracy by integrating multiple classifiers (Random Forest, Support Vector Machine, K-Nearest Neighbors) into a Voting ensemble method.
4. To evaluate and compare the performance of the proposed model against traditional classifiers, focusing on accuracy, false positives, and false negatives.

2. Literature Survey

Hassan Alizadeh, et al (2020) presented a novel GMM-dependent system for real-time network traffic classification and validation [11]. The models that were created were used to categorize the traffic and confirm where it was from.

Unlike traffic validation that primarily depended on probability testing to make sure that the application claiming itself as the traffic generator point was the actual source, the word classification meant the discovery of such an application through a maximum likelihood approach of highest posterior probability. With respect to traffic identification, the suggested way outperformed the current machine learning methods. The classification of the first 9 packets from a stream showed that the system was able to achieve a 97.7% accuracy. Also, when GMM was trained on just 0.5% of the streams, it achieved 0.966 accuracy.

Won-Ju Eom, et.al (2021) created a system for classifying internet traffic with the Software-Defined Networking (SDN) framework [12]. The designed system was completely put into action on the network manager using the controller's sharp computing capacity, joint discernibility, and programmability, which can provide exact, flexible, and practical traffic sorting. A range of universal criteria helped assessing the efficacy of the pertinent ensemble algorithmic strategies. Ensemble learning-based traffic classification techniques proved to perform better than others on real-world traffic datasets. One of the main differences between LVGBM and other deep learning models is the feature importance score margin and the time needed to train the model, which is significant. In particular, the use of the LightGBM classifier method yielded the best results.

Madhusoodhana Chari S., et.al (2019) evaluated the inevitability of separating protected web traffic in the modern scenarios. [13]. Sustainability in network management and raising traffic visibility were the ultimate objectives. Consequently, the early adopters of the traditional approaches had to face the practical challenges due to their high costs and lack of reliability and therefore, the only avenue that was explored in this area was the use of behavioral statistics to detect the flow of traffic. Anti-Overflow Network Layer data flow identifiers are defined based on security strength, bandwidth, and availability of minimal human intervention. In order to categorize traffic into numerous groups, a method that involves extracting a packet length signature was put forth. In order to determine the types of flowing traffic, this work used a feature set to train a novel J48 decision tree classification structure. Additionally, the framework's comprehensibility was examined.

Jing Ran, et.al (2018) developed 3D-CNN-based framework for traffic on networks classification [14]. This structure was an attempt to build on the developments in deep learning by applying the effective method of video analysis to network traffic classification. Effective classification could be achieved by the 3D-CNN by automating the attribute extraction step and identifying the best attributes following many validations. These features provide more thorough information than manually chosen traits. In order to solve the uncertainty of obtaining successful results via the most efficient feature extraction accompanied by the most successful classification approach in a low-collaboration setting, the work merged the data extractor and prediction model. Numerous experiments on the publicly accessible USTC-TFC2016 dataset verified that the proposed framework outperformed classic ConvNets with regard to accuracy. The main focus was on future studies on feature extractor integration with baseline benchmark categorization models.

Jiwon Yang, et.al (2019) provided a new payload-based classification technique that assessed the security of the transport layer by using unencrypted handshake packets sent back and forth between end hosts [15]. This method performed classification using a Bayesian Neural Network that was fed information about the TLS extension of the handshake packet, the compression technique, and the cipher suite. The advised procedure generated better results than other traditional payload-supported categorization frameworks in several tests. This method is to be applied in forthcoming research to sort out secure protocols that are there for further tests and observation.

Yu Wu, et.al (2018) suggested to use statistical multiplexing in MF systems in order to improve the conventional TDM-EPON framework and hence to increase the efficiency of upstream transmission [16]. They adopted two primary strategies in their works: (1) traffic classification, in which ML classifiers classified the upstream traffic into practical or impractical, and (2) idle-data transfer that managed to send the idle traffic so that no unnecessary EPON frames were sent. As a result, the first method, feature selection with classifier biases, investigated into two feature-selection strategies. In this respect, the one which was more effective was the improvement process through feature selection. The provided characteristics were used as input for the following methodology in the given approach. The suggested method did better than the others in terms of accuracy and resilience according to the experimental results of the dataset obtained. Thus, the proposed approach can reduce end-to-end latency to below 100 μ s and thus would establish user satisfaction compared to other alternatives in terms of traffic load and SNR.

Pratibha Khandait, et al (2020) developed a novel technology using the Deep Packet Inspection to categorize streaming traffic with a single payload scan [17]. As a matter of fact, the search was sublinear, implying a detailed heuristic was critical. The approach was heavily based on a dataset which contained streams from ten different apps in order to test it, which made clear that the new method is not only efficient but also precise. The JnetPcap library was utilized to execute the test in Java. In order to keep on improving even more, further work was intended to implement this concept in C leveraging the LibPcap package.

Guanglu Wei, et.al (2020) provided an architecture that involves Deep Learning (DL) using CNN in order to discriminate between network traffic on complex networks. A 2D grayscale image was produced for every traffic load of the network as the input [18]. This technique transformed network stream classification problems into image classification problems. This study advanced the study of using DNNs to classify network traffic. In comparison to traditional approaches, it offered the benefit of removing the requirement for extensive feature extraction and variable selection tasks, which enhanced prediction accuracy.

3. Research Methodology

This work's primary goal is to categorize the network traffic that is flowing. The following is an outline of the research methodology that encompasses different task to meet this goal.

1. Inputting Data set and Pre-processing: -The first step entails entering the dataset, which uses KDD—data collected from a trustworthy source. The 42-attribute NSL-KDD dataset was applied in this investigation. In order to improve the KDD'99 dataset and get rid of biased categorization outcomes, duplicate cases are eliminated. Only 20% of the training data is used, although many versions of the dataset exist. The data is represented as KDDTrain+_20Percent. To solve the issue of class imbalance, the SMOTE technique is applied.

2. Feature Extraction: - Establishing the connection between each of the attribute and the target set is the purpose of this phase. A false positive (FP) occurs when a sample appears normal but is identified as an intrusion, while a false negative (FN) refers to a situation where a case is actually an invasion but is identified as usual. If the system fails to detect an intrusion in the case of a false negative, it is considered a poor FN. Most Intrusion Detection Systems (IDSs) employ a layered approach, where if one layer misses a incident, subsequent layer attempts to do so. PCA is used at this stage to minimize features.

3. Classification: - Creating training and testing sets from the complete dataset is their primary objective of this step. The network traffic is categorized using a voting classifier. So as to effectively categorize the traffic, this technique integrates multiple classifiers. At this point, the test set is given the anticipated outcomes. Random Forest, SVM, and KNN are the classifiers utilized in this stage.

i. Random Forest: It is a kind of supervised machine learning (ML) method which is helpful for tackling the issue related to classify the data. While classifying the data, this algorithm emphasizes on dealing with the categorical data. To perform regression task, Random Forest (RF) aims at handling the continuous data. This algorithm is planned on the basis of an enormous amount of decision trees (DT), in which every DT is employed for predicting a class, and the class having a supreme number of predictions is considered as the predictive class of this algorithm. This algorithm is developed by following some particular stages in which every DT of this algorithm is consisted of a tree-like sequence of decision nodes. This demonstrated sequence is taken into consideration for partitioning the tree into a number of branches till the end node called leaf of the tree is attained. The final (leaf) nodes are utilized to helped in providing the predictive outcomes of every DT as output. In the end, the data is predicted by integrating outputs generated via multiple decision trees. This algorithm is effective to train the data at quick rate and avoid the issue regarding the overfitting. A technique known as bagging is adopted in this algorithm which enables every DT in selecting a sample dataset at random for attaining dissimilar kinds of trees at higher accuracy and least variance. Furthermore, the feature randomness (a technique in the ensemble) is put forward in order to maximize the divergence and mitigate the correlation amongst trees. This algorithm is expressed as:

$$RFfi_i = \frac{\sum_{j \in \text{alltrees}} \text{norm} fi_{ij}}{T} (1)$$

In which, $RFfi_i$ is used to define the significance of feature i . All the trees are exploited for computing this attribute, norm fi_{ij} is depicted with feature i whose normalization is done in a tree j , and T demonstrates the number of trees.

ii. SVM: SVM is a collection of supervised learning methods that use regression analysis to categorize data. In order for the learning process to assign a new categorical value as part of the prediction outcome, one of the variables in the training sample should be categorical. Since SVM ages the linear properties, it is a non-likelihood binary classifier lever. When used to high dimensionality, SVM is versatile and detects outliers in addition to classification and regression [1]. A training variable should ideally have at least two categories and be defined as follows:

$$x_i \in R^p, i = 1, \dots, n(2)$$

where x_i represents the training observation and R^p indicates the real-valued p -dimensional feature space and predictor vector space.

KNN: KNN is one of the simplest and oldest supervised machine learning algorithms used in classification; it classifies a given instance via the majority of the classes among its k -nearest neighbours found in the dataset. This algorithm relies on the distance metric used to determine the nearest neighbours of the given instance, and the most commonly used metric is the Euclidean distance, which is expressed in the following formula:

$$d(x_i, x_j) = \sqrt{\sum_{r=1}^n w_r (a_r(x_i) - a_r(x_j))^2} \quad (3)$$

where an example is defined as a vector $x = (a_1, a_2, a_3, \dots, a_n)$, n is the number of the example's attributes, a_r is r th attribute of the example and its weight is referred to as w_r , and $d = (x_i, x_j)$ are the two examples. To compute the class label of an example, the following formula is used:

$$y(d_i) = \underset{k^{x_j \in kNN}}{\operatorname{argmax}} \sum y(x_j, c_k) \quad (4)$$

where d_i is the example by which the algorithm will determine the class in which it belongs, the term x_j is one of the k -NNs present in dataset, and $y(x_j, c_k)$ indicates whether the x_j belongs to the class c_k . The result of Equation (2) is the class that has the most members of the k -NN, and is also the class wherein the example belongs. Euclidean distance is mostly used as a default distance in k -NN classification or k -means clustering to determine the " k closest points" of an example.

In this phase, these three classifiers Random Forest, SVM, and KNN are combined in the voting classifier approach to improve network traffic classification.

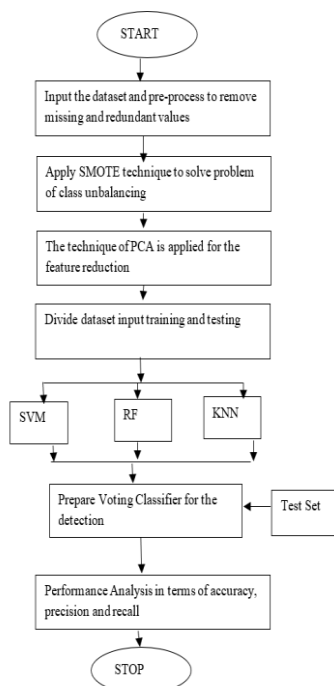


Figure 2: Proposed Methodology

4. Result and Discussion

The main focus of this study is network traffic classification. Data preparation, extracting features, and data classification are some of the stages involved in implementing the structure for classifying network traffic. The KDD dataset, which includes 42 attributes and a target set with numerous classes reflecting various attack kinds, is the dataset utilized for model testing. The competency level of the new approach is assessed based on several factors as discussed next.

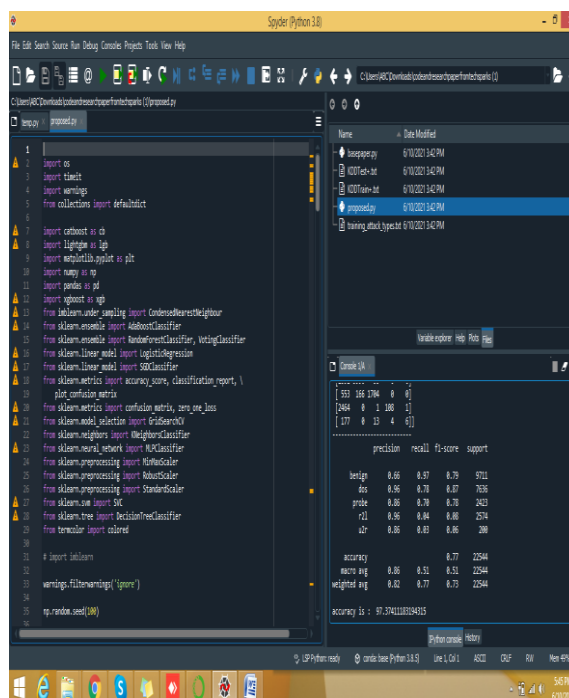


Figure 3: Execution of Presented Framework

The new approach is implemented with manifold methods, as seen in Figure 3. For the classification of network traffic, it combines SMOTE, PCA, and a Voting ensemble.

Table 1: Assessment of Performance

Model	Accuracy	Precision	Recall
SVM Classifier	75.74 Percent	81 Percent	76 Percent
Logistic Regression	72.67 Percent	80 Percent	77 Percent
KNN Classifier	70 Percent	72 Percent	76 Percent
Random Forest Classifier	75.78 Percent	76.89 Percent	75.90 Percent
Proposed Model	97.37 Percent	82 Percent	77 Percent

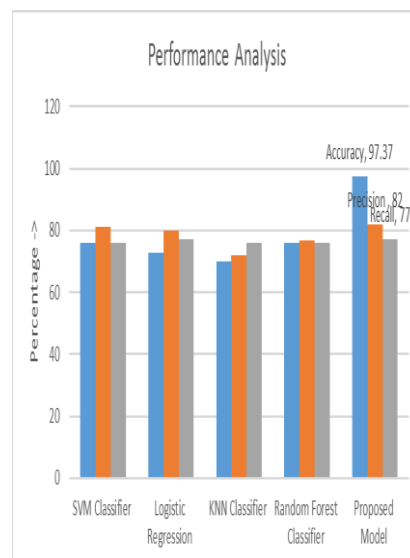


Figure 4: Performance Analysis

Figure 4 shows how the suggested method for classifying network traffic compares to a variety of classification models, such as SVM, LR, KNN, and RF. The suggested model exceeds the others with a maximum accuracy of 97.37%. This high accuracy shows how reliable the suggested framework is.

4.1 Comparative Evaluation

The new approach's performance is evaluated against standard classifiers like SVM, LR, KNN, and RF, showing a notable improvement in terms of accuracy and reduced false positives. The proposed model also demonstrates superior robustness in handling class imbalance and minimizing false negatives, as compared to traditional methods.

Conclusion

Network analysts improve data, heightening the challenge of detecting incidents. To develop a self-operating framework for identifying intrusions, effective and dynamic methodologies are required. These frameworks must be capable of learning and detecting evolving threats. Several extremely active, adaptable, and successful intrusion detection techniques are put forth for managing huge network traffic volumes. This paper tackles the issue of class misbalancing by introducing a novel architecture that integrates several strategies. The suggested model employs

PCA for feature reduction and the SMOTE technique to address the issue of class imbalance. Additionally, it incorporates several classifiers for classification, such as SVM, KNN, and Random Forest. This model attains about 97% accuracy, which is 15% better than previous models. Future research looks into the use of transfer learning for classifying internet traffic.

References

- [1] Jaehwa Park, JunSeong Kim, "A classification of network traffic status for various scale networks", 2013, The International Conference on Information Networking 2013 (ICOIN)
- [2] Ji-hye Kim, Sung-Ho Yoon, Myung-Sup Kim, "Study on traffic classification taxonomy for multilateral and hierarchical traffic classification", 2012, 14th Asia-Pacific Network Operations and Management Symposium (APNOMS)
- [3] Rui Yang, "The Comparison of Split-Flow Algorithms in Network Traffic Classification: Sequential Mode vs. Parallel Model", 2013, International Conference on Information Technology and Applications
- [4] ZebaAtique Shaikh, Dinesh G. Harkut, "A Novel Framework for Network Traffic Classification Using Unknown Flow Detection", 2015, Fifth International Conference on Communication Systems and Network Technologies
- [5] ShashikalaTapaswi, Arpit S. Gupta, "Flow-Based P2P Network Traffic Classification Using Machine Learning", 2013, International Conference on Cyber-Enabled Distributed Computing and Knowledge Discovery
- [6] Sung-Ho Lee, Jun-Sang Park, Sung-Ho Yoon, Myung-Sup Kim, "High performance payload signature-based Internet traffic classification system", 2015, 17th Asia-Pacific Network Operations and Management Symposium (APNOMS)
- [7] Yaojun Ding, "Imbalanced network traffic classification based on ensemble feature selection", 2016, IEEE International Conference on Signal Processing, Communications and Computing (ICSPCC)
- [8] Zhengwu Yuan, Chaozheng Wang, "An improved network traffic classification algorithm based on Hadoop decision tree", 2016, IEEE International Conference of Online Analysis and Computing Science (ICOACS)
- [9] Yang Hong, Changcheng Huang, BiswajitNandy, Nabil Seddigh, "Iterative-tuning support vector machine for network traffic classification", 2015, IFIP/IEEE International Symposium on Integrated Network Management (IM)
- [10] Chao Wang, Tongge Xu, Xi Qin, "Network Traffic Classification with Improved Random Forest", 2015, 11th International Conference on Computational Intelligence and Security (CIS)
- [11] Hassan Alizadeh, Harald Vranken, André Zúquete, Ali Miri, "Timely Classification and Verification of Network Traffic Using Gaussian Mixture Models", 2020, IEEE Access
- [12] Won-JuEom, Yeong-Jun Song, Chang-Hoon Park, Jeong-Keun Kim, Geon-Hwan Kim, You-Ze Cho, "Network Traffic Classification Using Ensemble Learning in Software-Defined Networks", 2021, International Conference on Artificial Intelligence in Information and Communication (ICAIC)
- [13] Madhusoodhana Chari S., Srinidhi H., Tamil EsaiSomu, "Network Traffic Classification by Packet Length Signature Extraction", 2019, IEEE International WIE Conference on Electrical and Computer Engineering (WIECON-ECE)
- [14] Jing Ran, Yexin Chen, Shulan Li, "THREE-DIMENSIONAL CONVOLUTIONAL NEURAL NETWORK BASED TRAFFIC CLASSIFICATION FOR WIRELESS COMMUNICATIONS", 2018, IEEE Global Conference on Signal and Information Processing (GlobalSIP)
- [15] Jiwon Yang, JargalsaikhanNarantuya, Hyuk Lim, "Bayesian Neural Network Based Encrypted Traffic Classification using Initial Handshake Packets", 2019, 49th Annual IEEE/IFIP International Conference on Dependable Systems and Networks – Supplemental Volume (DSN-S)
- [16] Yu Wu, Massimo Tornatore, Yongli Zhao, Biswanath Mukherjee, "Traffic classification and sifting to improve TDM-EPON fronthaul upstream efficiency", 2018, IEEE/OSA Journal of Optical Communications and Networking
- [17] PratibhaKhandait, NeminathHubballi, BodhisatwaMazumdar, "Efficient Keyword Matching for Deep Packet Inspection based Network Traffic Classification", 2020, International Conference on COMMunication Systems & NETWORKS (COMSNETS)

- [18] Guanglu Wei, “Deep Learning Model under Complex Network and its Application in Traffic Detection and Analysis”, 2020, IEEE 2nd International Conference on Civil Aviation Safety and Information Technology (ICCASIT)
- [19] Kumar, R., Singhal, N., & Chhabra, A. (2025). Revolutionizing Business Management Strategies for Enhanced Output Through the Integration of Deep Learning and Cloud Computing. *Journal of Information Systems Engineering and Management*, 10(58s).
- [20] Kumar, R., Singhal, N., & Chhabra, A. (2025). Hybrid Optimization algorithm with the combination of PSO and genetic algorithm for task scheduling in cloud computing. *E-Learning and Digital Media*, o(o). <https://doi.org/10.1177/20427530251331082>
- [21] Arpit Chhabra, Manav Bansal and Niraj Singhal, “Smart City-Shrewd Vehicle Versatility Utilizing IOT”, *International Journal of Engineering Trends and Technology*, Vol. 70, Issue 3, pp. 29-36, 2022.
- [22] Arpit Chhabra, Niraj Singhal and Syed Vilayat Ali Rizvi, “A Novel Algorithm of Safe-Route Traversal of Data for Designing the Secured Smart City Infrastructures”, *International Journal of Engineering Trends and Technology*, Vol. 71, Issue. 5, pp. 272-281, 2023.