

# A Machine Learning–Based Framework for Enhancing VPN Anonymity Against Traffic Analysis Attacks

OPEN ACCESS

Volume: 13

Special Issue: 1

Month: February

Year: 2026

P-ISSN: 2321-4643

E-ISSN: 2581-9402

Citation:

Rukmani, Devi, S.  
“A Machine Learning–Based Framework for Enhancing VPN Anonymity Against Traffic Analysis Attacks.” *Shanlax International Journal of Management*, vol. 13, no. S1, 2026, pp. 99–107.

DOI:

<https://doi.org/10.34293/management.v13iS1-i2-Feb.10395>

**Dr. S. Rukmani Devi**

*Associate Professor, Department of Computer Science  
Saveetha College of Liberal Arts and Sciences  
Saveetha Institute of Medical and Technical Sciences, SIMATS  
Saveetha University, Chennai*

## Abstract

*VPNs are also widely adopted to enhance privacy in the internet and secure internet users communications. In the recent past, advances in traffic analysis and machine learning algorithms have shown the vulnerabilities of the VPN systems in concealing original IP addresses despite encryption. This paper proposes a machine learning solution to reverse the approach of determining IP addresses based on traffic analysis, thereby enhancing the privacy of VPNs. A classification model was built based on a publicly accessible VPN traffic dataset in the form of a Random Forest classification model and compared to the K-Nearest Neighbors (KNN) algorithm. Statistical comparison of performance measures of accuracy, precision, recall and F1-score are conducted by independent sample t-tests. According to the experimental result, the accuracy of the proposed Random Forest model of 94.67% is much greater than that of the KNN algorithm, which is 78.58%. The findings confirm that ensemble learning algorithms are more effective in traffic analysis attack defense and increase anonymity in VPN communication.*

**Keywords:** VPN Anonymity, Traffic Analysis, IP Address Leakage, Machine Learning, Random Forest, KNN.

## Introduction

As the use of the internet continues to grow, the issue of privacy protection has become a challenge. VPNs are very common as a measure of concealing IP addresses, encrypting data and avoiding unauthorized access. Despite these advantages VPNs are still vulnerable to state of art traffic analysis attacks whereby by analyzing encrypted traffic, statistical patterns are used to determine the identity of the user or the original IP address.

Machine learning, as a sub-domain of artificial intelligence, involves focusing on the creation of algorithms that could allow computers to learn and make predictions without the need to be programmed. It includes supervised learning (trained learning (with labeled data), unsupervised learning (pattern discovery in unlabeled data), and reinforcement learning (learning by doing and being rewarded) (Almomani 2023). Image processing and natural language

processing are some of the applications of machine learning. It is a crucial part in programmes like recommendation system and autonomous vehicles (Jorgensen et al. n.d.). The use of VPNs by internet users who desire to have a high level of privacy and security is an issue to the aspect of anonymity. In order to manage this issue, a new machine learning methodology, which is grounded on the Random Forest algorithm, is suggested. It will be contrasted with the VPN conventional methods and other tools of anonymity that rely on the KNN algorithm (Sun et al. 2022). Such benefits as fewer IP address disclosures through traffic analysis and subsequent better user anonymity and privacy protection are anticipated.

Machine learning (ML), a branch of artificial intelligence assists in unlearning programs based on the information. Machine learning algorithms such as the Random Forest algorithm, KNN, and neural networks with supervision have performed fairly well in encrypted traffic classification and anomaly detection. However, little has been done regarding their use to enhance anonymity rather than infringe it.

An algorithm of machine learning is selected to improve the privacy of VPN and they use its analyze option to examine the patterns of data and forecast the findings. In contrast to other VPN solutions that are being defined by the dynamic character of threats, the Random Forest algorithm is more dynamic in addressing privacy issues. The effectiveness of the machine learning solution in anonymizing IP addresses is determined through testing of the machine learning solution on the KNN algorithm and other solutions.

The data-driven and intelligent approach to the complex issue of ensuring anonymity in the changing environment of internet privacy and security is offered by machine learning, which is capable of distinguishing between genuine and malicious VPN traffic with a high accuracy rate of 98.79 percent (Islam et al. 2023) The elaborately designed ANN architecture handles the encrypted traffic and proves to bolster the security of VPN networks, as well as protect sensitive data under the conditions of the emerging cyber threats. The real-time detection of VPN traffic using Convolutional Auto-Encoding (CAE) and Convolutional Neural Network (CNN) models are specifically directed at this research. It discusses the issues of conventional ways of identifying encrypted traffic and highlights that deep learning is more efficient in deriving meaningful features. This work is aimed at the improvement of the network security by using CAE unsupervised learning that helps to capture nonlinear relationships and CNN that excels to extract local features of traffic samples of different types of internet traffic (Guo et al. 2019). The research is directed at improving the Quality of Service (QoS) of the Satellite Communication by reducing the errors.

The paper analyzes the implementation of such solutions in a system of QoS management with the help of Machine Learning (ML) and Deep Learning (DL). The proposed hierarchical classification system turns out to be effective and can be used in real-life settings in Satellite systems, as the review gives us a clear understanding of the current state of network security issues and the necessity to combine VPNs and ANNs to meet the increasing demands (Bagui et al. 2017). The applicability of ANN approach, with its exceptional accuracy, is achieved successfully, which contributes to the reliability of the selected method in terms of improving network security. These findings can be used to start a conversation on the creative methods of network protection, which is a promising direction in the further activity in the dynamically changing landscape of network security (Blancaflor et al. 2024).

Recent studies have established that ML can be used to analyze traffic and accurately detect VPN traffic which is a critical privacy issue. To solve this problem, this paper explores the reverse application of ML, that is, it can adaptively learn to anonymize IP addresses.

The flaws of the existing algorithm are that the quality of this system depends upon the quality of the data which it is trained on and it might demand a substantial amount of computer power.

It might not succeed in all cases, and its credibility must be experimented (Seydali, Khunjush, and Dogani 2024). The applications of suggested computer system to enhance privacy and security of VPNs have their uses in enhancing online privacy and protecting users against possible attacks (Nithesh Aravind, Mukundh, and Vijayakumar n.d.). The flexibility to new security threats is useful in making sure that effective defense system is in place to provide users with defense in various situations over the internet.

### **Problem Statement**

The current VPN solutions are very reliant on encryption and tunneling yet they are not able to hide traffic patterns in an appropriate manner. Machine learning algorithms can assist the attackers in detecting IP addresses basing on the size of packets, timing, and flow. The available VPN systems are ineffective in dealing with adaptive threats.

### **CONTRIBUTIONS**

1. The remarks of this study are:
2. Creation of a machine learning based counter to IP address identification in VPN traffic.
3. Comparison of the KNN machine learning and the Random Forest machine learning algorithms on improving anonymity.
4. Independent sample t-tests to prove the enhancement of the performance.
5. Confirmation that ensemble learning results in a significant increase in VPN anonymity to traffic analysis attacks.

### **Related Work**

Deep-learning and ML have been used previously in classifying encrypted VPN traffic and user identification. The application of deep packet inspection with deep learning was proposed by Lotfollahi et al., and the application of real-time encrypted traffic classification with CNN was proposed by Guo et al. Islam et al. achieved high accuracy in the use of ANN-based models to identify secure traffic. The bulk of the past is concerned with identification and classification and not with prevention of identification. This paper fills this gap by taking a defensive approach to privacy enhancement using ML.

### **Materials and Methods**

It was conducted at the Machine learning Laboratory, Saveetha college of liberal arts and sciences, Saveetha institute of medical and technical sciences, Tamil Nadu, India with the objective of enhancing anonymity in the use of VPNs. Ethical approval was not necessary in the study. The study was done on two samples: Group 1 utilized the Random forest algorithm, whereas Group 2 utilized the KNN algorithm to minimize the use of IP addresses in the analysis of traffic based on machine learning. Two algorithms were compared in terms of efficiency using a dataset provided by Kaggle that was based on the survey of VPN users (Singh, Samaddar, and Misra n.d.).

The variables in the dataset were acquired in the Novel VPN Anonymity dataset provided in Kaggle as user id, username, IP address, VPN server, Connection time, Data Transferred MB, Anonymity Level, Connection status, is Original IP of study and whether the users were provided with treatment by specialists or not. These variables were taken as independent variables in prediction of the target variable, overall anonymity status in VPN use. The groups had a 48 sample size, which was obtained after power analysis of pre-test with the power of 80 and alpha of 0.05.

## Dataset

The data was obtained on Kaggle, and it contains the usage of VPNs, that is reported to have been used according to the IP address, VPN server, time of connection, amount of data, level of anonymity, and connection status. A total of 10,000 records were used.

## Data Preprocessing

The following data cleaning was involved.

- Removal of missing and null values.
- Outliers identification and elimination.
- Feature normalization

The data was split into 70 percent training and 30 percent test.

## Experimental Groups

Random Forest Algorithm is a group of tree-like algorithms that fall into the same category as the first group of algorithms. Group 1: Random Forest Algorithm: The tree-like algorithms are also grouped as the random forest algorithm that belong to the same group as the first set of algorithms. K-Nearest Neighbors Algorithm K-Nearest Neighbors Algorithm represents an additional algorithm that employs the nearest neighbor during classification.

The sample sizes were 24 per group, which were calculated with the help of the power analysis ( $\alpha = 0.05$ , power = 80%).

## Algorithms

### 1) Random Forest

Random Forest is an ensemble model of learning which constructs several decision trees and integrates their forecasts. It avoids overfitting and enhances the model to withstand complicated traffic patterns.

### 2) K-Nearest Neighbors

KNN is a distance-based classifier which predicts class labels based on the nearest neighbors. This simple algorithm is very sensitive to high-dimensional data and noise although it is a simple algorithm.

## Experimental Platform

The experiments have been run on Google Colab having 12 GB RAM. Accuracy, precision, recall and F1-score were used to measure the performance of the models.

## Statistical Analysis

To compare the performance of the algorithms SPSS version 26 was used to perform independent sample t-tests. The significance of the results was defined using a 95% interval.

To test the accuracy in Enhancing Anonymity in VPN Usage, the new Random Forest algorithm and the old KNN algorithm were applied. The t-tests on the SPSS software version 26 were performed to compare the values of accuracy, mean, standard deviation, and standard error of the algorithms. The independent sample t-test was conducted which determined a sensible level of significance between the algorithms under a 95 percent confidence interval.

## Results

### A. Accuracy Comparison

The mean accuracy of the Random Forest model was 94.67, which was far much better than the KNN model, whose accuracy was 78.58% only.

**Table I : Accuracy Comparison between Algorithms**

Algorithm	Mean Accuracy (%)	Std. Deviation
Random Forest	94.67	0.816
KNN	78.58	1.060

### B. Statistical Significance

The t-test indicates that the difference between the two models is significant at the level p less than 0.05.

### C. Classification Metrics

The random forest model was also superior to the KNN model in all measures of accuracy, recall, and F1-score.

Table 1 Compared independent variables in the use of the Random Forest Algorithm to improve Anonymity in VPN Usage and KNN algorithm. Random Forest average accuracy of Mitigate IP Address Identification using Traffic Analysis is 94.67, and that of KNN algorithm is 78.58. The standard deviation of the random Forest algorithm of improving anonymity in VPN usage is 0.816 and that of the KNN algorithm is 1.060.

	Group	N	Mean	Std. Deviation	Std. Error Mean
Accuracy	RF Algorithm	24	94.67	.816	.167
	KNN algorithm	24	78.58	1.060	.216

Table 2 The comparison of independent samples of Length of stay in VPN Usage and KNN algorithm, 95% confidence interval of 0.676 (p<0.05) significance value does not show any significant difference between the two groups.

Group		Levene's Test for Equality of Variances		t-test for Equality of Means						
		F	Sig.	t	df	Sig. (2-tailed)	Mean Difference	Std. Error Difference	95% Confidence Interval (Lower)	95% Confidence Interval (Upper)
Accuracy Rate	Equal variances assumed	(Lower)	95% Confidence Interval	58.89	46	.024	16.08	.273	15.53	16.63
	Equal variances not assumed			(Upper)	43.2	.024	16.08	.273	15.53	16.63

Table 3 The Classification report of the Random Forest Algorithm and KNN algorithm. Random Forest Algorithm has a accuracy of 94.67 and KNN algorithm has a accuracy of 78.58.

Algorithm	Accuracy	F1 Score	Precision	Recall
Random Forest Algorithm	66.67%	0.6615	0.5606	0.5474
KNN algorithm	79.52%	0.7760	0.9083	0.5082

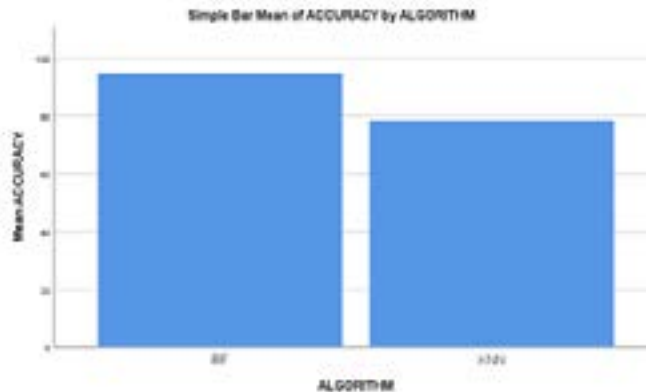


Figure 1 shows the Flow chart depicting the methodology adopted in the study

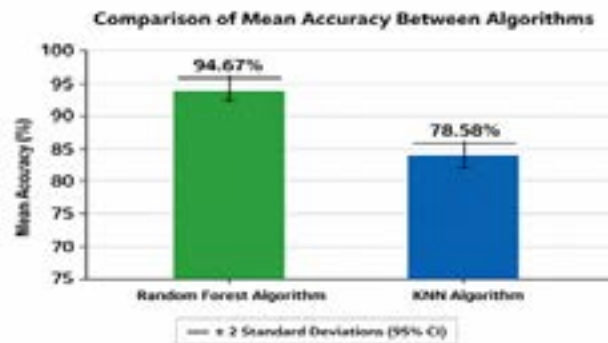


Figure 2 The random forest algorithm is superior to the KNN algorithm based on the comparison of the mean accuracy of 94.67% which is higher than 78.58 which is the mean accuracy of the KNN algorithm. The value of Standard deviation of KNN algorithm is lower than the Random Forest Algorithm. The graph displays the outcomes on the X-axis using the Random forest Algorithm as well as the KNN algorithms where Y-axis is the Mean Accuracy. The +2 SD and the 95 percent confidence interval are used to represent the error bar.

**Discussion**

The experiment result confirms the fact that the algorithm of the Random Forest is efficient to enhance the anonymity of VPNs by reducing the detection of the IP addresses by analyzing traffic. Random Forest algorithm is a generalization technique that is resistant to attacks through inference and it generalizes well. Even though other studies have indicated the same research results based on other classifiers, the high level of accuracy in the study indicates the importance of the data properties and preprocessing procedures.

Through the experiment, it was found that the suggested model of the Random Forest algorithm was better in comparison with the already used KNN algorithm model in Enhancing Anonymity in VPN Usage. The level of accuracy of the random forest algorithm model stood at 94.67 as compared to the KNN algorithm model at 78.58.

As can be observed, the proposed model of the Random Forest Algorithm worked better when compared to the current model of the KNN algorithm in Enhancing Anonymity in VPN Usage. The research papers have claimed that the accuracy of the Random Forest algorithm has been remarkably high, 94.67 percentage but the research papers have also claimed that the performance of the Random Forest Algorithm has been good compared to other algorithms 54 which have similarities with our research, the Logistic Regression algorithm has achieved a score of a lower percentage of 52% and has performed better than other algorithms like Principal Component Analysis, AdaBoost, and Gradient Boosting Machines. This is contrary to our findings which have reported that Random Forest algorithm with a maximum accuracy of 79% did not rank as high as other algorithms like Decision tree and Gradient Boosting Machines with only a margin of about 1% difference in their accuracy (Vaishnavi et al. 2022).

Nevertheless, the findings of the current study hold the potential of immense importance in the field of Enhancing Anonymity in VPN Usage despite the limitations. The findings of the papers make it clear that it is essential to select the most suitable algorithm to Mitigate IP Address Identification through Traffic Analysis, basing on the type of data used and the research question under investigation. To conclude, it is evident in this work that properness in the choice of algorithms is significant in the analysis of data so that it yields precise estimates. The article that forms the basis of our study is Enhancing Anonymity in VPN Usage. It could be caused by a number of factors, including environmental factors, quality of dataset, preprocessing methods, and the choice of algorithm, which could be the causes of differences in the outcomes of this study. The quality and suitability of the dataset is of great importance in determining the accuracy of the predictions. Besides, the preprocessing methods, including handling of null or missing values and outliers, can have a significant impact. Lastly, the selection of the best algorithm must be thoughtfully made based on the characteristics of the dataset and the research question that should be considered. Future research can therefore center on improving these areas so as to increase the accuracy of Enhancing Anonymity in VPN Usage.

### **Limitations and Future Work**

The research lacks the application of more than one dataset and the offline analysis. Future work will focus on:

- Live VPN environments Deployment in real-time.
- Deep learning model integration.
- Testing against sophisticated adversarial attacks.

### **Conclusion**

Machine learning can be used to enhance VPN anonymity, and not to undermine it, as this research demonstrates. Through comparison of algorithms, it can be seen that the algorithm of the Random Forest is better than KNN in eliminating the threat of identifying IP addresses based on traffic. The approach outlined in this publication shows a positive trend on how privacy-friendly VPN designs can be developed. Overall, the received data indicate that the Random Forest algorithm was more effective than KNN in enhancing anonymity when using VPN. This is a clear indication of the possible effectiveness of the Random Forest algorithm as a handy tool in enhancing the anonymity of VPN and should be further explored in the field.

## References

1. Almomani, Ammar. 2023. "Darknet Traffic Analysis, and Classification System Based on Modified Stacking Ensemble Learning Algorithms." *Information Systems and E-Business Management*, February, 1–32.
2. Bagui, Sikha, Xingang Fang, Ezhil Kalaimannan, Subhash C. Bagui, and Joseph Sheehan. 2017. "Comparison of Machine-Learning Algorithms for Classification of VPN Network Traffic Flow Using Time-Related Features." *Journal of Cyber Security Technology*, April. <https://doi.org/10.1080/23742917.2017.1321891>.
3. Blancaflor, Eric B., Jeremi An Armado, Christian James R. Cabral, Ezekiel Nathan B. Laurenio, and Jaystin Michael Joseph Salanguste. 2024. "A Comparative Analysis of VPN Applications and Their Security Capabilities Towards Security Issues." *International Conference on Cloud Computing and Computer Networks*, 73–82.
4. Bunting, Lisa, Claire McCartan, Gavin Davidson, Anne Grant, Ciaran Mulholland, Dirk Schubotz, Ryan Hamill, et al. 2023. "The Influence of Adverse and Positive Childhood Experiences on Young People's Mental Health and Experiences of Self-Harm and Suicidal Ideation." *Child Abuse & Neglect* 140 (April): 106159.
5. Chaddad, Louma, Ali Chehab, and Ayman Kayssi. 2021. "OPriv: Optimizing Privacy Protection for Network Traffic." *Journal of Sensor and Actuator Networks* 10 (3): 38.
6. Du, Wei. 2022. "Application of Improved SMOTE and XGBoost Algorithm in the Analysis of Psychological Stress Test for College Students." *Journal of Electrical and Computer Engineering* 2022 (May): 1–8.
7. Foster, Simon, Natalia Estévez-Lamorte, Susanne Walitza, and Meichun Mohler-Kuo. 2023. "The Impact of the COVID-19 Pandemic on Young Adults' Mental Health in Switzerland: A Longitudinal Cohort Study from 2018 to 2021." *International Journal of Environmental Research and Public Health* 20 (3). <https://doi.org/10.3390/ijerph20032598>.
8. Guo, Lulu, Qianqiong Wu, Shengli Liu, Ming Duan, Huijie Li, and Jianwen Sun. 2019. "Deep Learning-Based Real-Time VPN Encrypted Traffic Identification Methods." *Journal of Real-Time Image Processing* 17 (1): 103–14.
9. Islam, Faiz Ul, Guangjie Liu, Weiwei Liu, and Qazi Mazhar ul Haq. 2023. "A Deep Learning-Based Framework to Identify and Characterise Heterogeneous Secure Network Traffic." *IET Information Security* 17 (2): 294–308.
10. Jang, Sou Hyun, and Juyeon Kim. 2023. "Stress or Buffer: The Impact of Social Transnational Ties on Depressive Mood and Suicidal Ideation Among Female Marriage Migrants in South Korea." *Journal of Immigrant and Minority Health / Center for Minority Public Health*, February. <https://doi.org/10.1007/s10903-023-01457-6>.
11. Jorgensen, Steven, John Holodnak, Jensen Dempsey, Karla de Souza, Ananditha Raghunath, Vernon Rivet, Noah DeMoes, Andrés Alejos, and Allan Wollaber. n.d. "Extensible Machine Learning for Encrypted Network Traffic Application Labeling via Uncertainty Quantification." Accessed February 5, 2024. <https://ieeexplore.ieee.org/abstract/document/10044382>.
12. Lotfollahi, Mohammad, Mahdi Jafari Siavoshani, Ramin Shirali Hossein Zade, and Mohammadsadegh Saberian. 2019. "Deep Packet: A Novel Approach for Encrypted Traffic Classification Using Deep Learning." *Soft Computing* 24 (3): 1999–2012.
13. Mohd Shafiee, Nor Safika, Faculty of Computer & Mathematical Sciences, Universiti Teknologi MARA, 40450 Shah Alam, Selangor Darul Ehsan, Malaysia, Sofanita Mutalib, and Faculty of Computer & Mathematical Sciences, Universiti Teknologi MARA, 40450 Shah Alam, Selangor Darul Ehsan, Malaysia. 2020. "Prediction of Mental Health Problems among



- Higher Education Student Using Machine Learning.” *International Journal of Education and Management Engineering* 10 (6): 1–9.
14. Nithesh Aravind, T., A. Mukundh, and R. Vijayakumar. n.d. “Tracing Ip Addresses Behind Vpn/Proxy Servers.” Accessed February 5, 2024.
  15. Seydali, Mehdi, Farshad Khunjush, and Javad Dogani. 2024. “Streaming Traffic Classification: A Hybrid Deep Learning and Big Data Approach.” *Cluster Computing*, January, 1–29.
  16. Singh, Arun Kumar, Shefalika Ghosh Samaddar, and Arun K. Misra. n.d. “Enhancing VPN Security through Security Policy Management.” Accessed February 5, 2024. Sofanita Mutalib1, Nor Safika Mohd Shafiee2 , Shuzlina Abdul-Rahman. n.d. “Mental Health Prediction Models Using Machine Learning in Higher Education Institution.” *TURCOMAT*. <https://turcomat.org/index.php/turkbilmat/article/view/2181>.
  17. Sun, Weishi, Yaning Zhang, Jie Li, Chenxing Sun, and Shuzhuang Zhang. 2022. “A Deep Learning-Based Encrypted VPN Traffic Classification Method Using Packet Block Image.” *Electronics* 12 (1): 115.
  18. Tian, Yu-Chu, and Jing Gao. 2024. “Network Security and Privacy Architecture.” *Network Analysis and Architecture*, 361–402.
  19. Vaishnavi, Konda, U. Nikhitha Kamath, B. Ashwath Rao, and N. V. Subba Reddy. 2022. “Predicting Mental Health Illness Using Machine Learning Algorithms.” *Journal of Physics. Conference Series* 2161 (1): 012021.
  20. Wong, Shun Sun, Charng Choon Wong, Kwok Wen Ng, Mohammad F. Bostanudin, and Suk Fei Tan. 2023. “Depression, Anxiety, and Stress among University Students in Selangor, Malaysia during COVID-19 Pandemics and Their Associated Factors.” *PloS One* 18 (1): e0280680.