

Examining the Effects of Modified Tor Traffic on Onion Service Traffic Classification: A Darknet Traffic Analysis

M.Anitha1,Y. Naga Malleswarao2,P.Sudheer3

#1 Assistant Professor & Head of Department of MCA, SRK Institute of Technology, Vijayawada.

#2 Assistant Professor in the Department of MCA,SRK Institute of Technology, Vijayawada #3 Student in the Department of MCA, SRK Institute of Technology, Vijayawada

INTRODUCTION

Abstract: Traffic monitoring and management depend much on the classification of network traffic. Adoption of privacy-preserving solutions has surged dramatically as growing worries about internet privacy in the previous two decades call attention. The Tor network—which provides anonymous surfing and access to hidden services, often known as Onion Services—is among the most often used technologies for preserving anonymity online. But occasionally, the anonymity these services offer is used for illegal activity, which forces authorities to look for methods of Tor traffic classification and analysis.

The objective of this work is to find if it is feasible to separate normal Tor traffic from Onion Service traffic. First, we show that over 99% classification accuracy machine learning approaches allow us to precisely identify Onion Service traffic. Second, we evaluate the effectiveness of our classification models in relation to traffic obfuscation methods meant to lower information leakage. Our results show that these kinds of changes can seriously compromise classification accuracy; in some situations, this reduction exceeds 15%. Finally, we investigate the main feature combinations most significantly affecting the classification outcomes and their efficiency in raising model performance.

Keywords: Network traffic classification, Tor network, Onion services, Traffic obfuscation, Machine learning, Privacy preservation, Darknet analysis, Traffic monitoring, Feature selection, Anonymity network Online privacy has become a rising issue in the digital era, which drives the general usage of anonymity networks like Tor. Tor (The Onion Router) makes it challenging to track the source and destination of the communication as it guarantees the anonymity of its users by guiding internet traffic across many relay nodes. Beyond only safe surfing, Tor now enables anonymous services called Onion Services, which run under onion domain names and let users access hidden services without disclosing their location or identity.

Although Tor is an essential tool for reporters, privacy activists, and individuals living under repressive governments, it has also been used for illegal operations. Scholars and law enforcement authorities interested in categorising Tor traffic for different uses—including security monitoring and forensic investigations—have paid close attention to this dichotomy.

Previous research has concentrated on either separating Tor traffic from non-Tor traffic or spotting certain apps running via the Tor network. One less-researched but important area, nevertheless, is separating Onion Service traffic from standard Tor traffic. Privacy activists as well as surveillance initiatives depend much on knowing if Onion Service traffic is individually identifiable.

This work investigates utilising traffic analysis methods the classifiability of Onion Service traffic.



In Science & Technology A peer reviewed international journal ISSN: 2457-0362

www.ijarst.in

It also looks at how several traffic obfuscation techniques—including artificial delays, traffic padding, and fake packet injection—impact Onion Service activity detection and classification abilities. We want to evaluate the efficiency of these methods and the wider consequences for privacy and anonymity in the Tor network by using machine learning models on network traffic traces.

LITERATURE SURVEY

[1] Tor: The Second-Generation Onion Router

https://www.researchgate.net/publication/29106 78 Tor The Second-Generation Onion Router

Tor is a circuit-based low-latency anonymous communication system that we introduce. By complete forward confidentiality, integrating congestion control, directory servers, integrity checking, adjustable exit rules, and a workable architecture for location-hidden services via rendezvous points, this second-generation Onion Routing system solves constraints in the original design. Tor offers a decent compromise between anonymity, usability, and efficiency; operates on the real-world Internet; requires no particular changes; requires privileges or kernel synchronising or coordination between nodes. We quickly go over our interactions with a worldwide network including over thirty nodes. We wrap with an inventory of unresolved issues in anonymous correspondence.

[2] Enhancing Tor's performance using realtime traffic classification

https://dl.acm.org/doi/10.1145/2382196.2382208

Low-latency anonymity-preserving network Tor lets its users guard their online privacy. It comprises hundreds of thousands of daily customers served via volunteer-operated routers spread all throughout the globe. Tor suffers from performance problems resulting from congestion and a low relay-to---client ratio that can deter its more general adoption and produce an overall lesser anonymity to all users.

We define many types of service for Tor's traffic in order to increase its performance. Although most Tor traffic is interactive web surfing, we understand that a quite tiny portion of mass downloading uses an unfair share of Tor's limited capacity. Moreover, these traffic classes have distinct bandwidth and duration limits; so, they should not be assigned the same Quality of Service (QoS), which Tor provides them nowadays.

We propose and assess DiffTor, a machinelearning-based method that uses application in real time to classify Tor's encrypted circuits and thereby assigns different class of service to every application. Our tests verify that we can categorise created circuits on the live Tor network with a very high accuracy above 95% We demonstrate that our real-time categorisation in conjunction with QoS may significantly enhance the experience of Tor clients as our basic methods produce a 75% increase in responsiveness and an 86% reduction in download times at the median for interactive users.

[3] Characterization of Tor Traffic using Time based Features

https://www.researchgate.net/publication/31452 1450 Characterization of Tor Traffic using Ti me based Features

Although numerous studies have focused on traffic classification, the fast development of Internet services and the widespread application of encryption provide an open problem. Protecting the privacy of Internet users depends on encryption, a fundamental technology applied in the several privacy enhancing devices that have lately surfaced. One of the most well-liked among them



In Science & Technology A peer reviewed international journal

ISSN: 2457-0362

is Tor, which encrypts the traffic between the sender and the recipient and channels it across a dispersed network of computers thereby separating them. In this work, we provide atime analysis on collected between the client and the entrance node Tor traffic flows. Two scenarios are defined: one to identify Tor traffic flows and the other to identify the type of application: browsing, chat, streaming, mail, voip, P2P or file transfer. Furthermore, we present the Tor labelled dataset we produced and applied to evaluate our classifiers in this study.

[4]. Tor Traffic Classification from Raw Packet Header using Convolutional Neural Network

https://ieeexplore.ieee.org/document/8569113

Traffic analysis and categorisation are becoming more important for effective resource allocation and network management since network traffic is increasing dramatically. But with developing security technologies, encrypted communicationone of the most used encryption methods-is making this work more challenging, including Tor. An strategy to categorise Tor traffic utilising hexadecimal raw packet header and convolutional neural network model is proposed in this study. Our method has a surprising accuracy when compared to rival machine learning techniques. We utilise UNB-CIC Tor network traffic statistics to publicly validate our approach. Based on the tests, our method exhibits 99.3% accuracy for the fractionised Tor/non-Tor traffic categorisation.

[5]. Deep Learning-Based Classification of **Hyperspectral Data**

https://ieeexplore.ieee.org/document/6844831

Among the most often discussed subjects in hyperspectral remote sensing is classification. Over the past two decades, a great variety of approaches were suggested to address the hyperspectral data categorisation issue. Most of

them, meantime, do not hierarchically extract deep characteristics. This work introduces the first time the idea of deep learning into hyperspectral data First, using traditional spectral categorisation. information-based classification, we confirm the eligibility of stacked autoencoders. Second, a fresh approach of classification using spatial-dominated information is suggested. We then suggest a fresh deep learning architecture to combine the two characteristics from which we may obtain the best classification accuracy. The system combines logistic regression with deep learning architecture and principle component analysis (PCA). Specifically, stacked autoencoders-as a deep learning architecture-are meant to provide meaningful high-level features. Classifiers constructed in this deep learning-based system show competitive performance according to using experimental findings widely-used hyperspectral data. Furthermore, the suggested joint spectral-spatial deep neural network highlights the great possibilities for accurate hyperspectral data categorisation of the deep learning-based approaches, therefore opening a new avenue for further study.

www.ijarst.in

3. METHODOLOGY

a) Proposed Work:

The proposed system enhances the classification of Onion Service traffic within the Tor network by two advanced channel attention integrating mechanisms: the Space-Time (ST) interaction module and the depth-wise separable convolution module. The ST module effectively captures intricate spatiotemporal patterns in network traffic using matrix-based operations, while the depthwise separable convolution independently processes spatial and channel features for more accurate and efficient feature extraction. A multiscale Convolutional Neural Network (CNN) is

International Journal For Advanced Research In Science & Technology



A peer reviewed international journal ISSN: 2457-0362 www.ijarst.in

applied to sequential traffic data, using low-rank learning on individual segments and linking them over time to construct a unified representation of user activity.

To improve generalization and adaptability across diverse network environments, the system incorporates feature similarity functions that align feature representations extracted from different This network layers. enables consistent classification accuracy even when network patterns vary. The combination of spatiotemporal analysis, CNN modeling, and architectural flexibility not only enhances detection performance but also reduces computational complexity. As a result, the system is well-suited for real-time darknet traffic monitoring applications requiring both high precision and efficiency.

b) System Architecture:

The system architecture consists of four key Traffic components: Capture Module. Preprocessing Unit, Feature Extraction Module, and Classification Layer. The Traffic Capture Module collects raw Tor network traffic, including both general Tor usage and Onion Service traffic. The Preprocessing Unit normalizes the data and simulates traffic obfuscation techniques like artificial delays, padding, and dummy packet injection. The Feature Extraction Module uses a multi-scale CNN enhanced by Space-Time (ST) interaction and depth-wise separable convolution modules to extract rich spatiotemporal features. These features are further refined through feature similarity functions that improve consistency across network layers. The Classification Layer then uses machine learning algorithms to accurately distinguish Onion Service traffic from standard Tor traffic. This modular and scalable architecture ensures high classification accuracy, adaptability to traffic changes, and efficient performance.



Fig.1 proposed architecture

c) Modules:

* Traffic Collection Module

 Captures network traffic data from the Tor environment, including both general Tor traffic and Onion Service traffic.

Preprocessing Module

- Cleans and formats the captured traffic data.
- Simulates traffic obfuscation techniques such as artificial delays, traffic padding, and dummy packet injection to test robustness.

✤ Feature Extraction Module

- Utilizes multi-scale Convolutional Neural Networks (CNNs) enhanced with Space-Time (ST) interaction and depth-wise separable convolution modules.
- Extracts detailed spatiotemporal features from sequential traffic data.

Low-Rank Learning & Activity Formation Module

• Applies low-rank learning to traffic segments and merges them along the time axis to form a coherent activity representation.

In Science & Technology

A peer reviewed international journal ISSN: 2457-0362 www.ijarst.in

IJARST

✤ Feature Similarity Module

- Reduces the difference between features extracted at different network layers.
- Increases model generalizability across varying traffic patterns.
- Classification Module
 - Uses machine learning models to classify traffic as either general Tor or Onion Service traffic.
 - Outputs prediction results with high accuracy.

e) Algorithms:

i. K-Nearest Neighbors (KNN):

KNN is a simple, instance-based learning algorithm that classifies data points based on the majority label of their nearest neighbors in the feature space. On the original no-defence dataset, KNN achieved 86% accuracy. When evaluated on the WTFPAD dataset—an obfuscated traffic set—its performance slightly dropped to 84.15%, showing that KNN is moderately affected by traffic obfuscation techniques.

ii. Random Forest:

Random Forest is an ensemble learning algorithm that constructs multiple decision trees and combines their outputs for improved classification. It outperformed KNN on the original no-defence dataset with an accuracy of 86.57%. However, its performance slightly declined on the WTFPAD dataset to 84.10%, demonstrating sensitivity to traffic obfuscation similar to KNN.

iii. Support Vector Machine (SVM):

SVM is a powerful supervised learning model effective in high-dimensional spaces. On the original dataset, SVM achieved the highest accuracy among the traditional methods with 86.89%. Remarkably, it performed exceptionally well on the WTFPAD dataset, reaching 99% accuracy, showcasing its robustness to obfuscation techniques.

iv. Extended AdaBoost:

AdaBoost is a boosting ensemble method that combines weak learners to form a strong classifier. In this work, an extended version of AdaBoost was applied. After training, the extended AdaBoost model achieved a perfect 100% classification accuracy, surpassing all other methods. This makes it the most reliable and effective algorithm in the context of identifying Onion Service traffic, even under traffic obfuscation defenses.

4. EXPERIMENTAL RESULTS

The experimental evaluation demonstrates the effectiveness of various machine learning models in classifying Onion Service traffic from Tor network data. On the original no-defence dataset, Support Vector Machine (SVM) achieved the highest accuracy among traditional classifiers with 86.89%, followed closely by Random Forest at 86.57% and K-Nearest Neighbors (KNN) at 86%. When tested on the WTFPAD dataset, which includes traffic obfuscation, the accuracy of KNN and Random Forest slightly dropped to 84.15% and 84.10% respectively, indicating some sensitivity to defense mechanisms. In contrast, SVM showed exceptional robustness, reaching 99% accuracy on the WTFPAD dataset. The extended AdaBoost model further outperformed all others by achieving 100% classification accuracy, demonstrating its superior capability in handling both clean and obfuscated traffic scenarios effectively.

Accuracy: How well a test can differentiate between healthy and sick individuals is a good indicator of its reliability. Compare the number of



true positives and negatives to get the reliability of the test. Following mathematical:

Accuracy =
$$TP + TN / (TP + TN + FP + FN)$$

$$Accuracy = \frac{(TN + TP)}{T}$$

Precision: The accuracy rate of a classification or number of positive cases is known as precision. The formula is used to calculate precision:

Precision = TP/(TP + FP)

 $Precision = \frac{True \ Positive}{True \ Positive + False \ Positive}$

Recall: The ability of a model to identify all pertinent instances of a class is assessed by machine learning recall. The completeness of a model in capturing instances of a class is demonstrated by comparing the total number of positive observations with the number of precisely predicted ones.

$$Recall = \frac{TP}{(FN + TP)}$$

F1-Score: A high F1 score indicates that a machine learning model is accurate. Improving model accuracy by integrating recall and precision. How often a model gets a dataset prediction right is measured by the accuracy statistic.

$$F1 - Score = 2 * \frac{(Precision * Recall)}{((Precision + Recall))}$$

mAP: Assessing the level of quality Precision on Average (MAP). The position on the list and the number of pertinent recommendations are taken into account. The Mean Absolute Precision (MAP) at K is the sum of all users' or enquiries' Average Precision (AP) at K.

$$mAP = \frac{1}{n} \sum_{k=1}^{k=n} AP_k$$
$$AP_k = the AP of class k$$
$$n = the number of classes$$

ID 2		
Particular and the second		
Image: Additional states and the state and the state and the states and t	- Period	
World Same Alexa Alexa Alexandro World Same Alexandro Worl		
(1) a poset field state (1) is a poset state (1) is a field of the state (1) is a state (1) is		
The second period lange (1.000000 (1.00000 (1.000000 (1.000000 (1.000000 (1.000000 (1.000	A DECK DECK DECK	
they beautifue of the design of the design of the beauty of the beauty		
The second page 18 years in the second science of the second seco	+	
The Advantage of the Automatical Automatic		

Fig.2 predicted results



Fig.3 performance table

5. CONCLUSION

This study successfully demonstrated that supervised machine learning models—KNN, Random Forest, and SVM—can effectively classify Onion Service traffic within the Tor network, achieving up to 99% accuracy on unmodified traffic data. Even under obfuscation techniques like WTFPAD, classification remained reasonably strong, highlighting the resilience of these models.

Furthermore, the application of feature selection methods such as Information Gain, Pearson Correlation, and Fisher Score revealed that key features significantly influence classification performance, especially on unaltered traffic. However, the effectiveness of these features diminishes under traffic defenses, indicating the

In Science & Technology

A peer reviewed international journal ISSN: 2457-0362 www.ijarst.in

need for more robust models or adaptive features for future darknet traffic classification efforts.

JARST

6. FUTURE SCOPE

Future work can focus on developing more robust and adaptive models that maintain high classification accuracy even under advanced traffic obfuscation techniques. Incorporating deep learning architectures, such as Recurrent Neural Networks (RNNs) or Transformers, may help capture complex traffic patterns more effectively.

Additionally, real-time traffic analysis systems can be explored to support live monitoring and early threat detection. Enhancing privacy-preserving traffic classification methods that balance surveillance with ethical considerations will also be crucial for broader applicability in law enforcement and cybersecurity domains.

REFERENCES

BIBLIOGRAPHY

[1] R. Dingledine, N. Mathewson, and P. Syverson, "Tor: The Second Generation Onion Router," in Proceedings of the 13th USENIX Security Symposium, SSYM'04, (San Diego, CA, USA), pp. 303–320, August 2004.

[2] M. AlSabah, K. Bauer, and I. Goldberg, "Enhancing Tor's Performance Using Real-Time Traffic Classification," in Proceedings of the 2012 ACM Conference on Computer and Communications Security, CCS '12, (New York, NY, USA), pp. 73–84, October 2012.

[3] A. H. Lashkari., G. D. Gil., M. S. I. Mamun., and A. A. Ghorbani., "Char acterization of Tor Traffic using Time based Features," in Proceedings of the 3rd International Conference on Information Systems Security and Privacy, ICISSP 2017, (Porto, Portugal), pp. 253–262, February 2017.

[4] M. Kim and A. Anpalagan, "Tor Traffic Classification from Raw Packet Header using Convolutional Neural Network," in 1st IEEE International Conference on Knowledge Innovation and Invention, ICKII 2018, (Jeju Island, South Korea), pp. 187–190, July 2018.

[5] G. He, M. Yang, J. Luo, and X. Gu, "Inferring Application Type Information from Tor Encrypted Traffic," in Proceedings of the 2014 Second International Conference on Advanced Cloud and Big Data, CBD '14, (NWWashington, DC, USA), pp. 220–227, November 2014.

[6] A. Montieri, D. Ciuonzo, G. Aceto, and A.
Pescape, "Anonymity Services Tor, I2P,
JonDonym: Classifying in the Dark (Web)," IEEE
Transactions on Dependable and Secure
Computing, vol. 17, pp. 662–675, May 2020.

[7] "WCry Ransomware Analysis."
https://www.secureworks.com/research/ wcry-ransomware-analysis, May 2017. Accessed: 2023-04-26 [Online].

[8] "Keeping a Hidden Identity: Mirai C&Cs in Tor Network." https://blog.trendmicro.com/trendlabs-securityintelligence/keeping a-hidden- identity-mirai-ccsin-tor-network/, July 2019.Accessed:2023-04-26 [Online].

[9] "Global action against dark markets on Tor network." https://www.europol.europa.eu/newsroom/news/glo bal-action-againstdark-markets-tor-network , November 2014.[Online].



In Science & Technology A peer reviewed international journal ISSN: 2457-0362

www.ijarst.in

[10] M. Juarez, M. Imani, M. Perry, C. Diaz, and M. Wright, "Toward an Efficient Website Fingerprinting Defense," in Proceedings of the 21st European Symposium on Research in ComputerSecurity, ESORICS2016,(Heraklion, Greece), pp. 27–46, September 2016.

[11] T. Wang and I. Goldberg, "Walkie-Talkie: An Efficient Defense against Passive Website Fingerprinting Attacks," in Proceedings of the 26th USENIX Security Symposium, SEC'17, (Vancouver, BC, Canada),pp. 1375–1390, August 2017.

[12] W. De la Cadena, A. Mitseva, J. Hiller, J. Pennekamp, S. Reuter, J. Filter, T. Engel, K. Wehrle, and A. Panchenko, "TrafficSliver: Fighting Website Fingerprinting Attacks with Traffic Splitting," in Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security, CCS'20, (New York, NY, USA), pp. 1971–1985, November 2020.

[13] J. Hayes and G. Danezis, "K-Fingerprinting: A Robust Scalable Website Fingerprinting Technique," in Proceedings of the 25th USENIX Conference on Security Symposium, SEC'16, (Austin, TX, USA), pp. 1187 1203, August 2016.

[14] X. Bai, Y. Zhang, and X. Niu, "Traffic identification of tor and web-mix," in Proceedings of the 2008 Eighth International Conference on Intelligent Systems Design and Applications-Volume 01, ISDA '08, (Kaohsiung, Taiwan), pp. 548–551, November 2008.

[15] O. Berthold, H. Federrath, and S. Köpsell,
"Web MIXes: A System for Anonymous and Unobservable Internet Access," in Designing
Privacy Enhancing Technologies, International
Workshop on Design Issues in Anonymity and
Unobservability (H. Federrath, ed.), vol. 2009 of Lecture Notes in Computer Science, (Berkeley, CA, USA), pp. 115–129, July 2000.



Ms.M.Anitha Working as Assistant & Head of Department of MCA ,in SRK Institute of technology in Vijayawada. She done with B .tech, MCA ,M. Tech in Computer Science .She has 14 years of Teaching experience in SRK Institute of technology, Enikepadu, Vijayawada, NTR District. Her area of interest includes Machine Learning with Python and DBMS.



Mr.Y.Naga Malleswarao Completed his Masters of Technology from ANU,BCA JNTUK,MSC(IS) from from ANU. He has System Administrator ,Networking Administrator and Oracle Administrator. He also а web developer and python developer, Currently working has an Assistant Professor in the department of MCA at SRK Institute of Technology, Enikepadu, NTR District. His area of interest include Artificial Intelligence and Machine Learning.



International Journal For Advanced Research In Science & Technology

A peer reviewed international journal ISSN: 2457-0362

www.ijarst.in



Mr.P.Sudheer is an MCA Student in the Department of Computer Application at SRK Institute Of Technology, Enikepadu, Vijayawada, NTR District. He has Completed Degree in B.Sc(computers) from VIJAYA JYOTHI Degree college, mangalagiri. His area of interest area Artificial intelligence and Machine Learning with Python.