

DATA MINING AND DEEP LEARNING SYSTEMS FOR NETWORK TRAFFIC
CLASSIFICATION AND CHARACTERIZATION AT SCALE

by

Maya Kapoor

A dissertation submitted to the faculty of
The University of North Carolina at Charlotte
in partial fulfillment of the requirements
for the degree of Doctor of Philosophy in
Computing and Information Systems

Charlotte

2023

Approved by:

Dr. Siddharth Krishnan, Dr. Thomas
Moyer

Dr. Ahmed Helmy

Dr. Bei-Tseng Chu

Dr. Xiang Zhang

ABSTRACT

MAYA KAPOOR. Data Mining and Deep Learning Systems for Network Traffic Classification and Characterization at Scale. (Under the direction of DR. SIDDHARTH KRISHNAN and DR. THOMAS MOYER)

The real, complex network environment consists of an ever-increasingly diverse and large amount of data encapsulated in packets. Surveillance and monitoring of this traffic is a necessary task for law enforcement, cybersecurity, and intelligence agencies. Intercepted network traffic must be classified into multiple categories, such as the protocol encapsulation layers contained, application it originates from, user generating the traffic, and the traffic's malicious or benign nature. There is a lack of solutions which are able to classify packets individually without flow-based features. In order to address the gaps in current traffic classification and DPI techniques, we propose the initial release of the FORAGER toolkit, a software consisting of tools to extract hidden representations from individual packets and use these features in deep learning models to perform traffic classification. It uses data mining techniques to perform automatic generation of regular expression signatures, locality-sensitive hash fingerprints, and matrix and point cloud representations of packets. These are used as input features for corresponding deep learning models which can perform traffic classification on single packets in a real system. The models are multi-modal to capture multiple angles and dimensions of features for increased complexity of classification problems. They can be run in parallel for optimal throughput and scalability. Our experiments use these models in multiple configurations and scenarios to demonstrate superior performance and classification capability to advance the state of the art in complex network traffic surveillance and hidden representation learning.

TABLE OF CONTENTS

LIST OF TABLES	ix
LIST OF FIGURES	xi
LIST OF ABBREVIATIONS	1
CHAPTER 1: INTRODUCTION	1
1.1. Problem Statement	2
1.1.1. Data Complexity	2
1.1.2. Network Complexity	3
1.1.3. Scalability	5
1.2. Proposed Solution	6
1.3. Contributions	7
1.4. Summarization	9
CHAPTER 2: BACKGROUND	10
2.1. Deep Packet Inspection	11
2.1.1. Port-based Identification	11
2.1.2. Payload-based Identification	12
2.1.3. Flow-based Identification	13
2.2. Network Traffic Classification	14
CHAPTER 3: REXACTOR - REGULAR EXPRESSION SIGNATURE GENERATION FOR STATELESS PACKET INSPECTION	15
3.1. Introduction	15
3.2. Previous Work	17

3.3. Methodology	18
3.3.1. TRex: Apriori Tokenization of Packet Payloads	19
3.3.2. GRex: Genetic Sequence Alignment of Common Substrings	23
3.3.3. HRex: High Performance Regular Expression Scanning	27
3.4. Experiments and Results	27
3.4.1. Case Study A: HTTP Detection	30
3.4.2. Case Study B: VoIP Detection	31
3.5. Discussion and Limitations	32
CHAPTER 4: ALPINE/PALM - SIMILARITY SEARCH THROUGH LOCALITY-SENSITIVE HASH FINGERPRINTING	36
4.1. Introduction	36
4.2. Previous Work	38
4.3. Methodology	39
4.3.1. MinHash and Locality-Sensitive Hash Forest	39
4.3.2. ALPINE: A Locality-Based Packet Inspection Engine	41
4.3.3. PALM: Payload Analysis Using Locality Sensitive Measurements	41
4.4. Experiments and Results	42
4.4.1. Protocol Auto Detection	43
4.4.2. Traffic Type	43
4.4.3. TF-IDF Score Evaluation	45
4.4.4. Multi-Embeddings	46
4.4.5. Model Performance and Throughput	48

	vi
4.5. Discussion and Limitations	49
CHAPTER 5: MAPLE - IMAGE REPRESENTATION AND RECOGNITION OF INTERNET TRAFFIC	57
5.1. Introduction	57
5.2. Previous Work	59
5.3. Methodology	63
5.3.1. LeNet Model Configuration	64
5.3.2. ResNet Model Configuration	64
5.4. Experiments and Results	65
5.4.1. RTP Detection	66
5.4.2. Protocol Auto Detection	68
5.4.3. User Fingerprinting	69
5.4.4. Model Performance and Throughput	69
5.5. Discussion and Limitations	70
CHAPTER 6: DATE - POINT CLOUD REPRESENTATION AND DENSITY CLUSTERING ANALYSIS OF PACKETS	71
6.1. Introduction	71
6.2. Previous Work	71
6.3. Methodology	72
6.3.1. Point Cloud Creation	72
6.3.2. Density-Based Spatial Clustering of Applications With Noise (DBSCAN)	74
6.3.3. Packet Classification	76

6.4. Experiments and Results	77
6.4.1. RTP Detection	77
6.4.2. H.225 Detection	78
6.4.3. Model Performance and Throughput	78
6.5. Discussion and Limitations	79
CHAPTER 7: FORAGER - TOWARDS MULTI-EMBEDDINGS AND ENSEMBLE CLASSIFICATION	81
7.1. Introduction	81
7.2. Previous Work	83
7.2.1. VPN Traffic Profiling	83
7.2.2. Darknet Detection and Profiling	84
7.2.3. Filtering Compressed/Plaintext Traffic	85
7.2.4. Intrusion Detection over HTTPS	85
7.3. Methodology	86
7.3.1. Threat Model	86
7.3.2. System Design	90
7.3.3. Datasets	91
7.4. Experiments and Results	95
7.4.1. VPN Traffic Profiling	95
7.4.2. Tor Traffic Profiling	97
7.4.3. Plaintext Traffic over HTTPS	99
7.4.4. Compressed Traffic over HTTPS	100
7.4.5. SSH over HTTPS	103

	viii
7.4.6. DNS over HTTPS	103
7.4.7. Performance Metrics	104
7.5. Discussion and Limitations	105
CHAPTER 8: CONTRIBUTIONS AND DELIVERABLES	106
CHAPTER 9: CONCLUSIONS	108
9.1. Future Work	108
9.2. Summarization	109
REFERENCES	111

LIST OF TABLES

TABLE 3.1: Greek Symbols For GRex Encoding	26
TABLE 3.2: Supported Regular Expression Operators	26
TABLE 3.3: HTTP Signatures and Tokens	30
TABLE 3.4: HTTP Case Study Results for Regular Expressions	31
TABLE 3.5: HTTP Case Study Results for Tokens	31
TABLE 3.6: VoIP Signatures and Tokens	32
TABLE 3.7: VoIP Results for Regular Expressions	32
TABLE 3.8: VoIP Results for Tokens	32
TABLE 3.9: Worst-Case Space and Time Complexities for NFA and DFA [1]	35
TABLE 4.1: Time Complexities for LSHForest [2]	40
TABLE 4.2: Description of Combined Dataset	42
TABLE 4.3: Protocol Autodetection Results	44
TABLE 4.4: Multi-Hash Embedding Classification Results	47
TABLE 4.5: Performance Results for Varied Number of Classes	49
TABLE 4.6: Performance Results for ALPINE and PALM as a Multi-embedding Versus Hyperscan	50
TABLE 5.1: Configurations of each model used in the experiments	65
TABLE 5.2: Classification results of RTP vs non-RTP detection for ALPINE and PALM	66
TABLE 5.3: Classification results of RTP vs non-RTP detection for all MAPLE models	67
TABLE 5.4: Throughput results of RTP vs non-RTP detection for all MAPLE models	67

TABLE 5.5: Classification results of RTP vs non-RTP detection for MAPLE models A and C with additional training time	68
TABLE 5.6: Throughput results of RTP vs non-RTP detection for MAPLE models A and C with additional training time	68
TABLE 5.7: Results of multi-class detection for models A and C as macro averages	68
TABLE 6.1: Classification results of RTP vs non-RTP detection for ALPINE and PALM	78
TABLE 6.2: Classification results of RTP vs non-RTP detection for the DATE model	78
TABLE 6.3: Classification results of H.225 vs non-H.225 detection for ALPINE and PALM	79
TABLE 6.4: Classification results of H.225 vs non-H.225 detection for the DATE model	79
TABLE 7.1: ISCX VPN/non-VPN and Tor/non-Tor Datasets	93
TABLE 7.2: Summary of CIRA-CIC-DoHBrw-2020 Dataset	94
TABLE 7.3: Summary of eMews SSH over HTTPS Dataset	95
TABLE 7.4: Results for classifying VPN data	97
TABLE 7.5: Results for classifying Tor data by application	98
TABLE 7.6: Results for classifying Tor data by traffic type	98
TABLE 7.7: Results for VPN/Tor classification	100
TABLE 7.8: Results for encryption and compression detection	103
TABLE 7.9: Results for classifying DoH/non-DoH and SoH/non-SoH data	104

LIST OF FIGURES

FIGURE 3.1: Example SIP packet	17
FIGURE 3.2: High-level diagram of signature construction in REXACTOR. The system produces certain tokens and regular expressions which can be used with the HReX Scanner or any external filtering/sniffing application.	20
FIGURE 3.3: Modified Apriori Algorithm for string packet data using a sliding window algorithm. The sliding window process is repeated for all packet strings in the database file to make the full set of transactions.	22
FIGURE 3.4: Scoring matrix example of the Needleman-Wunsch Algo- rithm. Traceback follows from the right bottom corner to the top left cell. The alignment results in the maximum possible score = 0.	25
FIGURE 4.1: Jaccard Similarity of Shingle Sets	39
FIGURE 4.2: Example TCP/IPv4 header with extracted features for ALPINE outlined	41
FIGURE 4.3: Twenty-six Class Identification Problem using the ALPINE- PALM Combined Model	45
FIGURE 4.4: Confusion Matrix of Traffic Type Classification Results	46
FIGURE 4.5: Time and Space Complexity Results for ALPINE/PALM versus Hyperscan	51
FIGURE 5.1: Example application of kernel to pixelated packet data [3]	62
FIGURE 5.2: Grayscale images generated from packets in our combined dataset used in the following experiments	63
FIGURE 5.3: The architecture of our LeNet models	64
FIGURE 5.4: The architecture of our ResNet models	64
FIGURE 5.5: A confusion matrix of tracking user Bob across traffic streams in a changing network environment	69

FIGURE 6.1: 3D point cloud made out of a Session Initiation Protocol (SIP) packet	73
FIGURE 6.2: Illustration of types of points in DBSCAN	74
FIGURE 6.3: K-nearest neighbors graph generated for packets in the dataset using a <code>min_samples = 6</code>	76
FIGURE 7.1: A depiction of intelligence operations intercepting various types of traffic Bob is transmitting in connection to his extended criminal network	88
FIGURE 7.2: System diagram of FORAGER	90
FIGURE 7.3: Plaintext versus encrypted traffic results	101
FIGURE 7.4: Compressed versus encrypted traffic results	102
FIGURE 7.5: Compression type classification results	102

CHAPTER 1: INTRODUCTION

One of the most critical cybersecurity interfaces is the connection to the worldwide Internet. Over a petabyte of data is sent over the Internet every second. The majority of traffic is broken down into transportable, encapsulated units of data called packets. These packets contain multiple layers of transport and infrastructure information in headers which wrapper application and user data known as the payload. These layers are defined by Internet protocols, which are standardized by the Internet community in technical documentation. While payload data was originally assumed to be benign, malicious software downloads, viruses, and spyware are just a few examples of what may be hidden inside packet payloads. Attackers may also be trying to exfiltrate non-malicious but sensitive information such as enterprise or government secrets, national intelligence, or personally identifiable information. Furthermore, it is not just attackers surveying this data, as lawful interception must be performed for Internet regulation, network administration, counter-cyberterrorism, and criminal investigation.

In order to use deep learning and data mining techniques to perform network traffic classification, features must be extracted from packets to use as input and define classes. Techniques such as dynamic port allocation, encryption and compression of payloads, spoofing, and tunneling have significant impact on the ability of systems to extract meaningful features from packets. Furthermore, these features can be of high dimensionality or heterogenous and must be normalized and reduced to a processable, comparable format for machine learning algorithms. In this work, we explore techniques in data extraction and transformation of packets to uncover *hidden representations* for learning. We then develop and optimize models for processing and classifying data based on these derived features. The primary contribution of

this dissertation is the systemization of these techniques and models into a network traffic classification toolkit (**FORAGER** ¹) that has been released and made available for public use.

1.1 Problem Statement

The unsolved challenges in deep packet inspection that we address with FORAGER can be categorized into three sections: data complexity, network complexity, and scalability.

1.1.1 Data Complexity

Packet classification is not one-dimensional; the data is *complex* and *multi-faceted*. An individual packet may be accurately classified in multiple ways. For example, a packet may be sent by user Bob containing email data characterized by the POP3 application-layer protocol. Depending on the intent of the classifier, the same packet may be classified as email traffic, encrypted or compressed, benign traffic, traffic belonging to user Bob, or traffic containing POP3 data. Classification granularity could also be reduced to traffic type such as email or chat. This can be a more difficult problem as the features become more diverse across the individual classes. One limitation of existing research is that solutions tend to address only one classification problem or a subset of issues without recognizing or expanding to others. In contrast, we show in subsequent chapters how FORAGER considers the same datasets in multiple ways and can generalize to new and interesting queries.

An additional layer of complexity to packet data itself is modern encryption and compression of payloads. In order to properly process the data, one must first know if data is encrypted, compressed, or plaintext. This is not always obvious through header information or by port configuration, as it is possible to obfuscate these for malicious purposes. Analysts may find plaintext traffic on port 443, for example,

¹<https://pypi.org/project/forager-toolkit/>

where they usually expect encrypted traffic. Similarly, they may be able to decompress traffic if they can distinguish it from similarly entropic encrypted traffic. Being able to filter out and then further investigate this content could reveal important information related to surveillance or forensic operations.

Finally, payloads or protocol standards themselves are not always strongly signatured. Particularly data protocols such as FTP (file transfer protocol) Data or RTP (realtime transfer protocol) can be highly entropic and therefore hard to identify even with state-of-the-art, regular expression scanning or by indicators in the packet header. FORAGER uses hidden representation learning to better understand and distinguish this data through generative features.

1.1.2 Network Complexity

Identifying traffic through a single packet can be difficult also because of network complexity. Protocol and data exchanges in the modern Internet are complex over the lifetime of the session exchange, particularly for large amounts of transferable data. Routines like Voice-over-IP (VoIP) calls, file transfers and downloads, streaming services, and peer-to-peer sharing can transmit hundreds of thousands of packets of highly variable data over several minutes to hours. These exchanges are also not guaranteed to be one-to-one; for example, content servers may be queried by hundreds of thousands of users at once for mass streaming services.

Deep packet inspection for the purpose of identifying traffic type or protocol is often performed for sessionization or reconstruction of whole data flows (for example, reconstructing a video call). Individual packets from multiple protocols such as SIP (session initiation protocol) and RTP must be correctly identified and associated together in real time to capture and correctly decode the full message.

Another complication for DPI is obfuscated network or routing information due to network architecture's advances in anonimity. VPNs (virtual private networks) are private network infrastructures configured on top of public Internet services. VPNs

can be used for increased security, privacy, and anonymity. Furthermore, administrators may control how data is transmitted and what data can be accessed by whom according to the network configurations. Depending on the site configuration, layers of tunneling, or dynamic mapping of networks, VPNs can add layers of Internet protocol (IP) encapsulation which obfuscates addresses. *IP-over-IP encapsulation* means that the initial protocol stack of a stream can contain several layers of IP which may change over time. In many intrusion detection systems, surveillance equipment, and other middlebox technologies, 3-tuple or 5-tuple information made up of source and destination addresses and port numbers are used to track streams. If layers are added or changed, this can interrupt stream tracking and cause unintended breaks or loss of streams. Attackers or covert operators may wish to avoid tracking by switching tunnels, as well. Without the ability to track content streams regardless of network layer changes, the data law enforcement and the intelligence community care about will become lost and untraceable.

Tor ("The Onion Routing") [4] is an internet framework intended to combat fingerprinting to prevent third-party tracking and surveillance through multiple layers of encryption. The Tor browser can be freely downloaded and used by anyone to achieve virtual network anonymity and access to sites which may be otherwise blocked for certain users or regions. It accomplishes multiple layers of encryption through a series of network hops where each server provides unique encryption keys which must be wrapped on and then peeled off at each stop. Over two million people employ Tor services as of 2022; The vast majority of Tor usage is benign but studies show a significant portion of users do perform illegal activities on the network and disproportionately so in countries which have less Internet regulation and government surveillance [5]. Thus, it can be useful for surveillance and forensic operations to identify Tor traffic and necessary for DPI systems employed to do so.

There are also active adversaries who work to confuse DPI and avoid tracking.

In order to bypass firewalls or avoid detection, threat actors may choose to send traffic over non-standard ports or attempt to spoof traffic as a particular content type by sending it over a known port. There are many reasons to send non SSL/TLS-encrypted traffic over port 443. One primary motivation is that most firewalls default-allow traffic on ports 443 and 80. Thus, traffic may pass through alongside other data without the need for additional configuration or permissions. Devices which may wish to configure with a local network with minimal user or administrator overhead, such as Internet of Things (IoT) devices, have been discovered to send cleartext traffic alongside the encrypted traffic on these ports. Where security is not regulated or enforced, encryption standards also may just be ignored for simplicity's sake. Because the traffic is typically encrypted, payload-based deep packet inspection is often left unemployed on port 443 in favor of using compute resources on other traffic channels. Thus, port 443 has the potential to be used as a covert data channel if left unscanned. This has been realized in the wild; research has uncovered applications of foreign origin on public university networks running traffic through port 443 in order to avoid firewall detection. Threat actors may route data intended for other commonly scanned ports, such as email traffic and file transfer data, through 443 to avoid content-based inspection. Tunneling protocols like SSH may also be run using HTTPS to create encrypted covert data channels. In this work, we discuss this exact threat scenario and employ FORAGER to profile traffic on port 443, providing a clearer picture to nefarious activity which may be occurring.

1.1.3 Scalability

Machine learning and deep learning based systems are plagued by computational complexity and stream buffering requirements. In much of the existing research, entire traffic flows or certain portions of flows are required before classification can begin. In a system at scale processing terabits of data per second, it is not practical to buffer this many streams in dynamic memory. Furthermore, classification speeds and models

must be capable of accelerated computation through hardware or simple enough to be run in parallel quickly to keep up with line rate. For real-time classification, it may not be possible to wait for multiple packets in a single stream to identify a particular flow. Even in offline forensic analysis, entire streams may not be available or recoverable. Rather, the system must make a best effort guess without knowing the state of the traffic as to what application layer protocol is being carried for immediate parsing and processing. If possible, classifying the application layer traffic per packet would provide the lowest latency and highest throughput possible in the system.

1.2 Proposed Solution

The process of network traffic classification can be broken down into steps: deep packet inspection, representation learning, and machine learning-based classification. The task of identifying or classifying network traffic can begin with shallow or deep packet inspection. In these steps encapsulation layer and payload information of the packet is analyzed for matching content or features indicative of a particular class. This data may be further expanded through forms of representation learning or data transformation, where hidden features may be uncovered and used as input to the next step, machine learning. These algorithms such as neural networks or state vector machines learn these representations as classes and use the embedded data for the future testing phase in order to appropriately group incoming data. At runtime, a result for classification is returned on the input data. In this work, we provide novel contribution in applying several new forms of representation learning through data mining and transformation steps, as well as appropriate deep learning classifiers following the creation of the hidden features.

This work proposes FORAGER as a toolbox for network traffic classification using data mining and deep learning methods. It is a culmination of lessons learned through the development of these modules and appropriately addresses the problems discussed in the scope of this dissertation. For scalability, FORAGER uses single-packet inspec-

tion techniques which may also be run in parallel for ultimate optimization. We demonstrate high success and negligible performance impact in protocol autodetection, traffic type and application profiling, and user fingerprinting. The system is also able to adapt to port obfuscation, spoofing, and tunneling architecture and identify VPN and Tor routing. Other types of data information such as compression, compression type, and encryption are able to be profiled through the FORAGER models.

The scope of this work focuses on traditional network data. Specialized datasets such as Internet of Things networks, mobile networks, and vehicular networks are not considered here but would be future opportunities for expansion. We also intentionally do not consider traffic flows or flow-based features. The data inputs and classifications performed in this work are done on a per-packet basis for minimized latency and buffering. Hardware optimization and acceleration has been proposed in other literature but are outside this scope. The focus of this work is the software implementations of data engineering, transformation, and analysis. Scalability is assessed through runtime metrics such as testing speed, throughput of data through the system (transformation plus classification), and dynamic memory usage. The models are also gauged with performance metrics from machine learning such as precision, recall, and F1-score.

1.3 Contributions

As part of this dissertation, PCAP sources containing twenty-six different application layer protocols were processed by our **TAPCAP** solution and combined into a diverse collection of tabularized data. We provide this dataset as a public resource for the community ². As a deliverable, we also publish FORAGER 1.0 ³ along with support packages for REXACTOR 1.0 ⁴ and TAPCAP 1.0 ⁵. The goal of providing

²<https://github.com/mayakapoor/protocol-dataset>

³<https://pypi.org/project/forager-toolkit/>

⁴<https://pypi.org/project/rexactor/>

⁵<https://pypi.org/project/tapcap/>

these software packages and corresponding user documentation is to encourage the research community to use these tools in their own experiments in order to further apply what we have learned in work accomplished in this research. Each of the data engineering tools and deep learning tools provide intellectual contribution to the field in the following ways.

REXACTOR is our automatic regular expression generation solution, furthering the state-of-the-art by encoding substrings through global sequence alignment and mining frequent tokens in an original algorithm for more enriched, expressive regular expressions than previous work. **PALM** (Payload Analysis using Locality-Sensitive Measurements) and **ALPINE** (A Locality-Sensitive Packet Inspection Engine) are innovative methods for creating locality-sensitive hash embeddings from network traffic which is a unique alternative to regular expression scanning previously used for deep packet inspection. Compared to the previous method, the tree-based search scales sublinearly and the fixed hash representations require a linear amount of storage space. These methods accurately identify many classes such as traffic type, protocol type, application, and more. Our Matrix PayLoad Encoder, **MAPLE**, and Density-based Analysis Tensor Encoder, **DATE**, explore applications of image-based and density analysis-based embeddings to an unsolved problem of RTP detection and profiling data protocols. The DATE model is a unique approach to three-dimensional, point cloud modeling of packets which is not covered in previous literature and proves to have applications in our experiments.

Each model is assessed with measures of precision, recall, and F1-score in different classification problems. We also provide comparison between models in order to better understand which transformations work better for different network environments and different traffic types or problems. In several of the experiments we also compare against state-of-the-art work such as the Hyperscan regular expression scanning library or other machine learning profiling techniques and systems to show our system

contribution. Our system provides an innovation also in its multi-modality, which we demonstrate in a case study where we analyze traffic in multiple ways on port 443. For scalability, we also assess memory usage and throughput across scenarios to propose the system at scale so that users may be able to design their classification strategy to their own requirements and available resources.

1.4 Summarization

In the following chapters, we first explain the background work behind deep packet inspection and the path toward deep-learning based traffic classification. Chapter 3 introduces our work into automatic regular expression signature generation (REXACTOR), which learns commonalities between packet payloads, transforms them into regular expressions, and performs text-based automata searching of the data. In chapter 4, systems based on locality-sensitive hashing (ALPINE and PALM) introduce data compression and fixed storage space along with a decision-tree based forest classifier. We also introduce multimodality, a recurring enhancement to classification in our system. chapters 5 and 6 introduce MAPLE and DATE respectively, which perform two-dimensional and three-dimensional transformations of data. In the MAPLE architecture, we further apply CNN classification. In DATE, density-based clustering analysis is performed and statistics used as input to a neural network classifier. Finally, the ensemble system introduced in chapter 7 is the culmination of the previous models and the basis for FORAGER as a useable application for real cyber engineers and surveillance specialists. In chapter 8 of this work, we present the FORAGER software and documentation resources. Chapter 9 concludes with further research opportunities and a summarization of issues discussed.

CHAPTER 2: BACKGROUND

This section will introduce previous work related to deep packet inspection and the progression toward machine learning for network traffic classification. Related work toward specific techniques will be further discussed in the “Previous Work” subsections of subsequent chapters.

As a disclaimer to this section and a warning to researchers exposed to other works in this problem area, we emphasize that not all network traffic classification systems are comparable to one another. Systems can classify the same data using different input features or into different classes. Choorod et al [6] and architects of DIDarknet [7] classify the same ISCX Tor/non-Tor dataset in different ways. Choorod et al focuses solely on darknet detection, whereas DIDarknet further explores traffic profiling by application and traffic type. This dataset along with the ISCX VPN/non-VPN dataset are used by many works, but different systems require certain input features. Not all systems are applicable to all network environments due to classification capability and system requirements such as processing speed or memory usage. Depending on the implementation, a trade-off to prioritize accuracy over resource utilization may not always be the right choice (or vice versa). It may be essential to prioritize accuracy if the system is capable of storing and processing the amount of data at an acceptable rate. Having multiple systems capable of using different features is a more practical approach to the network traffic classification problem because it is unlikely one solution will work equally as well for multiple problems, and many real world applications may not afford all the features that one system requires over another.

2.1 Deep Packet Inspection

Deep packet inspection is the process of analyzing traffic data as it comes across the network. Packet headers contain information which can tell what type of application layer data is being transported in the packet. Furthermore, packet payload contents may be inspected for signatures matching protocols, applications, users, or other content related to a particular class. More recently, systems have explored data transformation such as embeddings or statistical analysis on raw packet data in order to generate hidden features for representation learning-based DPI.

2.1.1 Port-based Identification

IANA, the Internet Assign Number Authority, requires default ports for certain protocols in standard traffic [8]. Classifying traffic based on port number is increasingly unreliable due to non-standard port usage and dynamic port translation. Applications or users may ignore protocol standards for both legitimate and illegal reasons. For example, a network administrator may choose to load-balance traffic to a different port for resource allocation. On the other hand, criminal activities could be conducted on a non-standard port for obfuscation purposes. Sometimes, benign IoT traffic is sent over a port like 443 or 80, as well [9]. Some applications are also known for not following standards. For example, Skype is not HTTP traffic but uses Web TCP to listen on port 80 [10]. A web designer can also host HTTP content on port 8080 instead of 80. Peer-to-peer networks may also conduct their operations on any range of ports including those intended for other types of traffic [11].

Some protocols or traffic sessions may dynamically configure their ports from a given range; this also poses a problem to port-based identification in a single-packet context. Network Address Port Translation (NAPT) allows multiple computers in a subnet to use one global IP address between the original source and final destination. The port translation occurs when a sender is listening over a specific port that is not

available within the NAPT set. In this scenario, the port is translated to an unused number and then mapped into a translation table. This way, the destination port will match what the NAPT set on the destination side, not the port specified originally by the sender [12].

2.1.2 Payload-based Identification

Because port-based identification became inaccurate in real world applications, researchers moved toward inspecting other parts of the packet, particularly the payload. In shallow packet inspection, we focus on what information can be extracted from the physical, data link, network, and transport layers as described by the OSI model [13]. This is then extended into deep packet inspection which analyzes the session, presentation, and application layers of a packet.

DPI can be rendered ineffective by encryption protocols like TLS, but is still used today for plaintext traffic [14]. While much of the Internet’s traffic today is encrypted, there are still applications for plaintext analysis afforded by deep packet inspection. Control, setup, and session management messages may be sent in the clear and may provide contextual information or protocol type indication [15, 16, 17]. Public Wi-Fi access points can also pass unencrypted traffic [18]. Internet of Things (IoT) devices have rapidly permeated our society and introduced new protocols and formats, yet their growth currently outpaces the security regulations for them. Popular market devices like Google Home and Amazon Echo are so rapidly flooding the market that they are not yet regulated by cybersecurity laws or encryption standards and privacy policies put in place. Thus, they often send and store data messages in plaintext [19, 20, 21]. Factory SCADA and PLC networks and older legacy systems may also not support encryption [22] and thus may be inspected in the clear.

Even in encrypted payloads, researchers have found that signatures can still exist. For example, Skype signaling messages are end-to-end encrypted but still contain a specific 15-byte sequence at the beginning of a Skype login server message [23].

HeaderHunter [24] uses packet size, direction, and TCP information in order to perform DPI pattern matching. This and other systems [25, 26] are not single-packet classification solutions, and rather require a full session capture in order to classify packets. The focus of this work and strength of our solution is the ability to classify packets per-packet, without the need for flow-based features.

2.1.3 Flow-based Identification

While we focus on packet-based features, there is an entire field of literature and research devoted to flow-based identification which has yielded impressive results. There are other heuristics from protocol flows which developers and researchers may derive further information or hidden features from. These statistics-based identification methods typically focus on mean, minimum, and maximum of packet sizes, the number of packets seen in a given exchange, burst rates, the time between packets, TCP flags, IP/port numbers, flow duration, and so forth [25, 14, 27]. In the existing literature, most use machine-learning or statistics techniques that study time and flow-based features of packet streams [28, 29, 30, 31, 32, 33, 34]. Some mechanisms require flow-based or time-series features, which means capturing and analyzing several packets in a flow [35, 36], or in some cases the entire session [37, 38, 39, 40, 41, 42, 43, 44]. The need to capture or buffer entire flows for classification does not scale well to large systems with memory or CPU constraints. Furthermore, it delays classification in real-time applications which may be trying to dynamically process the data. There is also no guarantee to receive the whole stream if capture begins in the middle of the session or there is data corruption. In our work as in some others, we focus on per-packet classification. In this scenario, we propose solutions which require only a single packet at any point in the flow for classification. We do not feed forward any information to the next classification or record any statistics about the packets over time or between arrivals. While this reduces the number of possible features and in some cases accuracy on public datasets, this is a constraint of the problem scenario

we are considering in this work.

2.2 Network Traffic Classification

Common substrings [45] or features like packet statistics can be extracted and used to train and make informed decisions in machine learning-based models. Quality of Service requirements and adaptive streaming features have been extracted to classify video streams using a Naive Bayes algorithm [46] or state vector machine (SVM) [31, 29]. These potentially have significant drawbacks such as high computational complexity, high-retraining time, and reduced performance compared to other models like K-Nearest Neighbors [28]. C4.5 and other decision tree methods have proven performant when used to classify traffic on large datasets, as well [47, 33]. Stateless applications of neural networks have also been able to achieve greater accuracy than the C4.5 approach [38]. Similarly, word embeddings have also become practical in the stateless context with Packet2Vec [48].

CHAPTER 3: REXACTOR - REGULAR EXPRESSION SIGNATURE GENERATION FOR STATELESS PACKET INSPECTION

3.1 Introduction

Regular expressions, often referred to as regexes, are sequences of characters used to specify search patterns in text for string-based searching algorithms. Stephen Cole Kleene first developed regular expressions as a derivative of study in automata theory and formal languages in theoretical computer science. Today, regular expressions are powerful notations used in many tools including UNIX programs like `grep`, `awk`, `sed`, and `vi`; search engines like Google, Yahoo, and Bing, and deep packet inspection software like Snort, RE2, and TRE. Universal standards exist for defining regular expressions, such as Perl, PCRE, and POSIX syntaxes.

Regexes specify a subset of possible strings derived from the alphabet based on conditions described by the syntax of the expression. Operators exist which define the patterns. For example, a vertical bar (`gray|grey`) would indicate alternatives and match the string literals “gray” or “grey.” The regular expression `gr(a|e)y` would be considered an equivalent regular expression and also match both these literals. Quantifiers specify how many occurrences of a character will appear in the given string. The question mark indicates zero or one occurrences of the preceding character. For example, `(too?)` would match “to” or “too.” The Kleene star specifies zero or more occurrences of the preceding character; the regex `a*` would thus match a string containing zero up to an infinite number of *as*, versus the Kleene plus regex `a+` which would match strings containing one up to an infinite number of *as*. A specific number of occurrences or a range may also be specified with curly brackets. These are just a few examples of regular expression operators, not including POSIX groups or other

such syntax-specific patterns.

Regular expression pattern matching is a text-based search solution which works best with data which is strongly signed. Specifically, the method names, keywords, common return codes, and string tokens that frequently occur in packets make it possible for network engineers to come up with regexes that can match specific protocols for application layer identification. One example of strongly signed packet data which can be matched with regular expressions is HTTP headers. The regular expression, `(GET|HEAD|POST). * HTTP/1.(1|0)x0Dx0A`, will match all HTTP packets containing a method signature like “GET/ index.html HTTP/1.0”, or “HEAD /index.html HTTP/1.1.”

Machine learning research has found that for protocols which have commonalities across messages, unique features tend to be found in the first N bytes of packet payload [49]. In many mainstream application protocols such as POP3, SMTP, HTTP, SIP, XMPP, IRC, and others, the first line of payload data contains methods, commands, version numbers, protocol names, and other potential features. Figure 3.1 shows this starting line in the SIP protocol. Various protocols call this feature by different names; for example, HTTP calls this line “request-line” or “response-line” [50]. FTP refers to this line as “request” or “response” command [51]. For consistency, we refer to this portion across protocols as the starting line. In testing other protocols which do not necessarily follow a request-response paradigm, we adapted the framework to include other message type options. For example, XMPP utilizes three kinds of stanzas which present unique features: IQ, message, and presence [52].

Creating regular expressions for content filtering or pattern matching is a complex problem that often requires subject matter expertise and manual intervention. Cross-examining multiple packet payloads and examining technical reference documents like RFCs is a time-consuming but necessary process to derive accurate and effective regular expressions. Furthermore, these signatures must be regularly maintained, tested,

Starting Line	INVITE sip:123-456-7890@sip.com SIP/2.0
SIP Header	Via: SIP/2.0/UDP 10.0.1;branch=a2dG4wL1234 Max-Forwards: 50 To: <sip:098-765-4321@sip.com> From: <sip:123-456-7890@sip.com>;tag=123456 Call-ID: 68ce48d219a CSeq: 1 INVITE Contact: sip:10.2.0.1:5060 Content-Type: application/sdp Content-Length: 100

Figure 3.1: Example SIP packet

and updated. Regular expressions become immediately ineffective for example if version numbers are added or changed. In the cybersecurity domain, malware signatures may be easily subverted through simple code or string manipulation so permutations must also be considered. In order to expedite this process and increase signature scope, researchers have leveraged knowledge discovery techniques and applied data mining to extract common features from packet payloads and header contents as input into regular expression signature generation algorithms.

We designed REXACTOR, a **R**egular **EX**pression **A**priori **C**onstruc**TOR** to extend the state-of-the-art and create more expressive signatures than other solutions are previously capable of generating. REXACTOR uses techniques from bioinformatics and natural language processing to encode signatures which use more regex characters and capture more data than previous work while reducing the dynamic memory footprint when employed in regular expression scanning solutions.

3.2 Previous Work

Applications for regular expression signature matching to packet payloads originate in content inspection for malware detection. Early work in automatic signature generation for polymorphic worms such as Autograph [53] provided a foundation for signature-based analysis. This expanded into detecting personally identifiable information as well as policy violations such as copyright infringement, inappropriate or

sensitive information on enterprise networks, and censorship. Additionally, Tang et al [54] found multiple sequence alignment to be effective in rewarding consecutive substring extractions and tolerating noise in traffic. SigBox [45] introduced the technique of generating substring tokens from packet payloads and applying Apriori data mining to find the most frequent ones. The related CSP Algorithm [55] and works by Wang et al [56], Szabo et al [57], LASER [58], and AutoSig [59] all include feature extraction and sequence alignment. Wang et al describe their system as a four-stage process; we generalize this concept to all these solutions and our own in order to compare various approaches.

In the first step, data is pre-processed in order to prepare it for use at the next stage. This typically involves extraction of certain fields or N number of bytes from the packet payload. For all these systems, sessionization of packet flows and sometimes defragmentation and reassembly based on TCP/IP header values is required [45, 56, 55, 57, 59, 58]. Once the data is collected and prepared, the second stage finds common substrings across packet data and/or protocol flows. The systems examined which perform feature extraction use either a subtree approach [56], a longest common substring algorithm [59, 58], motifs [57], or sequential pattern mining [45, 55]. Third, a method of alignment is used to align data based on commonalities between packets. Wang et al, Szabo et al, Vinoth et al [60] and LASER use bioinformatics approaches to perform these alignments. A substring tree can also be used for this purpose [59]. Some of these alignment strategies are additionally informed by a scoring matrix influenced by the tokens derived in the previous step [56, 57, 58]. Once the sequences are aligned in an optimal manner, in the fourth and final stage the systems convert their results into regular expressions.

3.3 Methodology

We use a modified version of Apriori algorithm combined with frequency position tables in order to calculate single byte-length tokens and frequently occurring sub-

strings using choice operators and position constraints [56]. REXACTOR also uses genetic sequence alignment to find commonalities in payloads and encodes this alignment in order to create a regular expression from it. Our system performs better than previous work in memory space allocation by additionally using position frequencies to specify character classes and length constraints instead of generic wildcards. This creates real-time performance optimization in the scanner. Furthermore, unlike previous solutions our system requires no packet buffering or traffic flow information to do its identification.

REXACTOR is made up of component modules which perform various stages of data processing and analysis in order to create optimal regular expressions. The flow of the system is shown in Figure 3.2. REXACTOR takes in as input packet capture (PCAP) files and extracts session payload strings for supported protocols. The input can be of any mix of traffic types, and is therefore well-suited to learning in the wild from real network captures. This extraction layer takes in parameters from the command line interface to select a supported protocol and message type to specifically isolate training data from the mixed input. It also offers an extraction option for the full data layer for analysis. For our experiments, we selected HTTP, SIP, and RTSP as protocols and extracted request and response session payloads for each of them.

3.3.1 TRex: Apriori Tokenization of Packet Payloads

Inside the starting lines, some substrings appear frequently across packets and thus can be used as string literals in signatures. In the SIP protocol, the line “SIP/2.0” appears at the end of all starting lines in all our training data SIP requests. In HTTP, the protocol name and version “HTTP/1.1” appear consistently, as well. We use a modified version of the Apriori Algorithm [61] in order to find frequently occurring substrings, or *tokens*, in our packet strings. Shim et al also use Apriori in SigBox for content strings [45].

Apriori is an algorithm used for finding frequent item sets and association rules

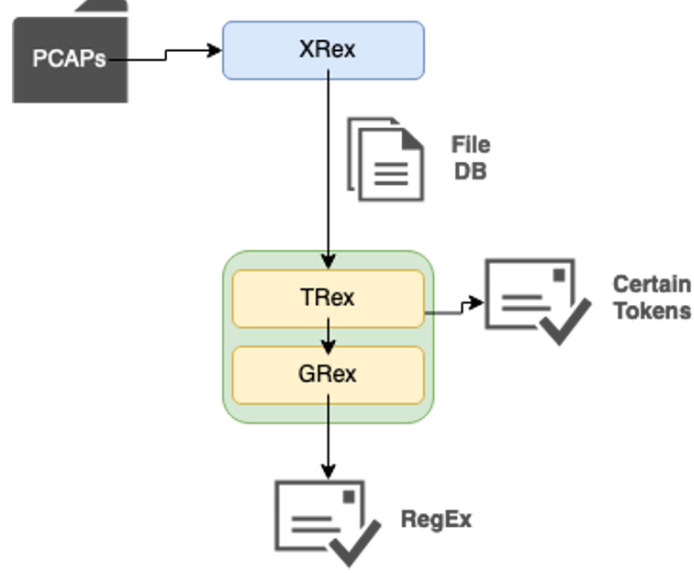


Figure 3.2: High-level diagram of signature construction in REXACTOR. The system produces certain tokens and regular expressions which can be used with the HReX Scanner or any external filtering/sniffing application.

from databases. The algorithm uses a bottom-up approach to combine items into incrementally larger item sets, keeping only those sets which meet some minimum threshold of support. The support of an itemset in the database is defined as the number of instances of the itemset in the transactions over the total number of transactions [61].

$$\text{Support}(X) = \frac{\text{Instances of } X \text{ in transactions}}{\text{Number of transactions}}$$

In the formal definition provided by Rakesh Agrawal [62], let $I = \{i_1, i_2, i_3 \dots i_n\}$ be a set of items of length n and $D = \{t_1, t_2, t_3 \dots t_n\}$ be the set of transactions known as a database. Every transaction t_i has a unique identifier in D and $\forall t_i : x \in t_i \rightarrow x \in I$, or every transaction is a subset of items in I . Apriori relies on the anti-monotonicity of the support of itemsets which assumes that all subsets of a frequent itemset must

also be frequent, and all supersets of an infrequent itemset must also be infrequent.

Algorithm 1: Modified Apriori Algorithm [62]

Result: Frequent Itemsets
 $k = 1, I_k = \{\text{frequent item sets of size } 1\};$
while ($I_k \neq \emptyset$) **do**
 $C_{k+1} = \text{apriori_gen}(I_k);$
 for $t \in D$ **do**
 $C_t = \text{subset}(C_{k+1}, t);$
 for $c \in C_t$ **do**
 $c.\text{count}++;$
 end
 $I_{k+1} = \{c \in C_{k+1} \mid c.\text{count} \geq \text{minsup}\};$
 end
 $k = k + 1;$
end
tokens = $\cup_k I_k;$
return $\text{prune_substrings}(\text{tokens});$

We mine tokens using a modified Apriori approach. Each packet string derived from pre-processing is treated as an individual line in the database. TRex uses a sliding window algorithm of length k to create substring items which are grouped into a list treated as a single transaction. We initialize the window size to $k = 2$ in order to prevent likely frequent but meaningless single-byte tokens and set the maximum itemset length $l = 1$. The first iteration calculates the frequency of substrings of length $k = 2$ and each subsequent loop increments k until no more frequent itemsets (substrings) are found.

Because the order of substrings matters in regular expression matching, we limit the maximum itemset length to 1. We increase item size instead by increasing the window size, thereby preserving order of characters in the modified algorithm. At each iteration, we also prune the resulting token set by replacing tokens i in the current set $S = \{I_{k_0}, I_{k_1}, I_{k_2} \dots I_{k_n}\}$ with any strings j in the result set $I_{k_{n+1}}$ which are superstrings of i and $\text{support}(i) = \text{support}(j)$. This reduces redundancy of tokens such as “SIP” and “SIP/” when the sets are unioned together.

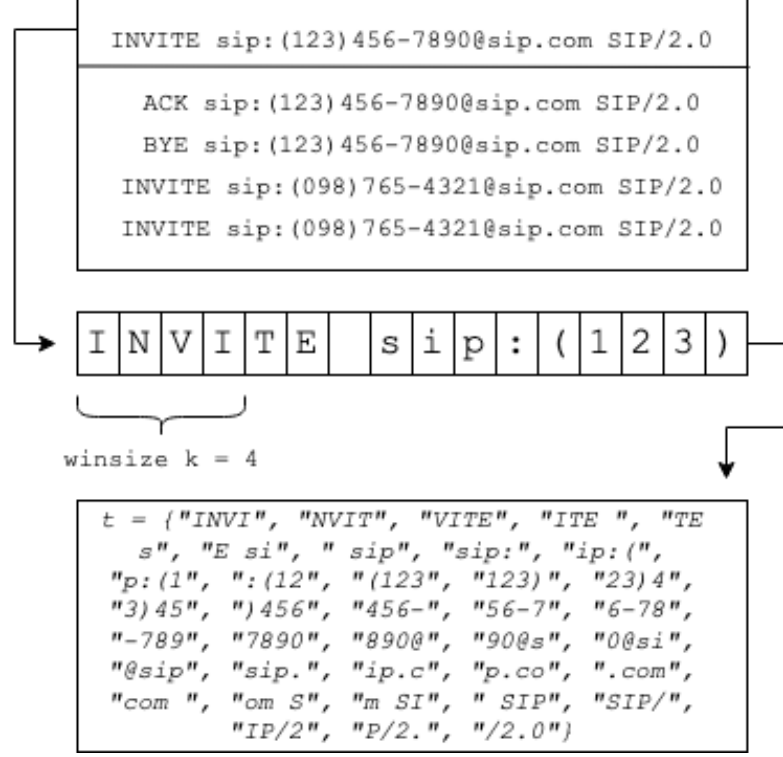


Figure 3.3: Modified Apriori Algorithm for string packet data using a sliding window algorithm. The sliding window process is repeated for all packet strings in the database file to make the full set of transactions.

Certain tokens. Substrings with support equal to 1.0.

Frequent tokens. Substrings with support equal to a manually configured threshold. For example, in SIP 2.0 the method words INVITE, ACK, and BYE frequently occur in request packets.

In order to give positional context to tokens, we introduce position frequency tables from natural language processing (NLP). This allows TRex to compute prefix and suffix groups from frequent tokens. Given a list of tokens $t_1, t_2, t_3, \dots, t_n$, TRex constructs a position frequency table for each token in each starting line in the database. The resulting vector $V = [\{t_1 : pos_1, t_2 : pos_2, t_3 : pos_3, \dots, t_n : pos_n\}, \{t_1 : pos_1, t_2 : pos_2, t_3 : pos_3, \dots, t_n : pos_n\} \dots]$ is referenced for tokens which occur at $pos = 0$. If the sum of support for those tokens found at $pos = 0$ is approximately 1.0 with respect to a minimal noise threshold, all tokens in this subset are captured in a prefix group. In the SIP

request example, the methods `INVITE`, `ACK`, and `BYE` would be captured and separated with regular expression choice operators and a beginning position constraint in the following manner: `^(INVITE|ACK|BYE)`.

For suffixes, TRex uses string manipulation to reverse both the tokens and starting lines, which intuitively reverses the position significance also so that $pos = len(starting_line)$ is now decides what goes into the capture group. Frequent tokens are reversed back to their original order before being inserted into the suffix, and choice operators and an end position constraint are appended. Example suffixes may include response codes, status codes or messages, or version numbers of protocols. An example output from TRex configured for HTTP requests and run on data containing both v. 1.0 and v. 1.1 requests may derive a suffix such as `(HTTP/1.0|HTTP/1.1)$`.

In addition to forming prefixes and suffixes, position frequency tables also give meaning to single-byte tokens which occur with $freq = 1.0$ at a fixed position. Wang et al found that some protocols have distinctive features at fixed offsets. As an example, they showed that the QQ messaging application always ends its packets with the byte `0x03`, indicating the end of a message [56]. We calculate position frequency tables similar to their approach for each character in each starting line and repeat the process in reverse in order to find single-byte tokens. These tokens are then added to the certain tokens set for future regular expression insertion. The substring tokens are then extracted from the relevant database lines and the remainder is preserved for the alignment stage.

3.3.2 GRex: Genetic Sequence Alignment of Common Substrings

GRex uses principles from genetic sequencing algorithms in order to pairwise align the remaining starting line data. Because we pre-process data into starting lines and TRex removes the found prefixes and suffixes to concentrate on only data which should be aligned, it is appropriate in our system to use a global alignment algorithm which attempts to align the entire sequence pair. We also use progressive alignment

to sequentially align pairwise sequences to derive a resulting optimal alignment.

We chose to utilize the Needleman-Wunsch Algorithm for scoring and global alignment [63]. As a dynamic programming algorithm, Needleman-Wunsch utilizes a 2D scoring matrix which rewards matching characters and penalizes mismatches according to predefined parameters. Furthermore, a gap penalty is introduced to account for insertions/deletions (indels) in the string alignment. The principle of the algorithm is to greedily maximize the score of alignment per cell. Each cell in the matrix also contains a pointer value which points to the origin of the highest score and will aid in the final trace-back stage [64]. The score function F with gap score g , sequence $A = \{a_1, a_2, a_3, \dots, a_n\}$, and sequence $B = \{b_1, b_2, b_3, \dots, b_k\}$ is defined as follows:

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + \text{score}(a_i, b_j) \\ F(i-1, j) + g \\ F(i, j-1) + g \end{cases} \quad (3.1)$$

In the initialization phase, the first row and column are set as the gap score times the distance from the origin, which is the upper left most corner of the matrix. From here, the algorithm recursively computes the match score, horizontal gap score, and vertical gap score for each cell in the matrix. The match score is calculated from the preceding diagonal cell's score and the score of alignment of the two characters (match or mismatch). The horizontal gap score is the sum of the score to the left and the gap score. Similarly, the vertical gap score is the sum of the cell above and the gap score. The cell is set to the maximum of these values and the pointer to the cell (left, vertical, or diagonal) where the maximum came from. Once the matrix has been filled, the pointers are used to trace back through to find the highest-scoring alignment [64]. Figure 3.4 demonstrates this scoring algorithm in a matrix.

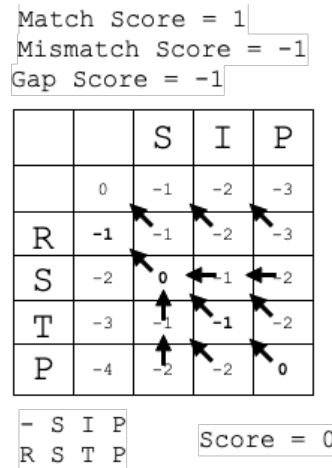


Figure 3.4: Scoring matrix example of the Needleman-Wunsch Algorithm. Traceback follows from the right bottom corner to the top left cell. The alignment results in the maximum possible score = 0.

GRex modifies the Needleman-Wunsch algorithm to encode Greek symbols into our alignments which represent specific character classes, insertions and deletions, and mismatches which will later be decoded in the final translation step and used to generate regular expressions. Table 3.1 provides a key to GRex’s Greek symbol encoding. Previous work used an encoding algorithm for indels and mismatches [56], but we created our own encoding schema for GRex and expanded this to also incorporate defining character classes. Specifically, we encode symbols based on whether a character in a given position is or is not a white space character, alphabet character, or digit. Table 3.2 shows the regular expression operators GRex derives from these encodings at the translation step.

The encoding process runs in a progressive manner, such that two strings are aligned, compared, and encoded, and that return is then passed again and compared against another string in the sequence. Once aligned using the Needleman-Wunsch algorithm, the character types of each index are compared. Depending on the type of character, the returned string from this process is appended with a representation of either the character itself if the indexes are matching or a Greek character encoding

Table 3.1: Greek Symbols For GRex Encoding

Φ	insertion/deletion
Δ	mismatch
ψ	gap
Σ	alphabet character
ω	not alphabet character
Π	digit
λ	not digit
Ω	white space character
σ	not white space character

Table 3.2: Supported Regular Expression Operators

Type	Symbol	Definition
character class	.	matches any character
	[]	matches any character in brackets
	[^]	matches any character not in brackets
	\s	white space character
	\d	digit
	\w	word character
	\S	not white space character
	\D	not digit
	\W	not word character
quantifier	*	matches zero or more times
	?	matches at most one time
	+	matches one or more times
	{ <i>m</i> , <i>n</i> }	matches between <i>m</i> and <i>n</i> times
	{ <i>m</i> , }	matches <i>m</i> or more times
	{, <i>n</i> }	matches at most <i>n</i> times
other	^	matches beginning of data
	\$	matches end of data
	()	capture group
		OR operator

otherwise. Once the new string is fully appended, it is returned to the loop. This returned string is compared to the next entry within the data. This process loops until every line from the packet data is aligned into the final encoding.

In order to create the regular expression, the replacement function is run against the generated string. Each index is checked and replaced with the regular expression counterpart with correct special character escaping. The final regular expression is constructed from the symbolic string. Frequent tokens are combined with choice operators and positional constraints \wedge and $\$$ for prefixes and suffixes, respectively. Prefixes are appended to the beginning of the generated regular expression, and suffixes to the end.

3.3.3 HReX: High Performance Regular Expression Scanning

The final component is a sniffing application written in C++ which reads in PCAP files using `libpcap` and performs regular expression scanning with Intel’s Hyperscan library [65]. Hyperscan provides significant performance improvement over other commercial but popular tools like Snort due to its reduction of duplicate operations in string and pattern matching and use of single-instruction, multiple data (SIMD) operations. In their experimentation, creators of Hyperscan found it outperformed Snort by a factor of 8.7 [66]. Our string literal tokens created by REXACTOR’s TRex can be applied in a pre-filter system like Snort, as well. Furthermore, Hyperscan does additionally support string literal, or keyword, matching. We provide both solutions for flexibility in real-world applications as the network capture environment is rarely one-size-fits-all. As our tokens are learned via unsupervised learning, they do not require manual analysis and thus scale well to larger, unknown data sets.

3.4 Experiments and Results

We performed two case studies using REXACTOR regular expressions and tokens. The first is a comparison-based analysis using our own HTTP signature and ones

developed by related works. Additionally, we assessed REXACTOR’s ability to develop a signature and keywords for SIP and RTSP, two protocols found in VoIP traffic flows. We use measures of precision, recall, and F1 score in order to assess signature efficacy.

We use measures of precision, recall, and F1 score in order to assess signature efficacy. Our measure of recall is the number of packets n correctly identified as protocol type P divided by the sum of the number of packets n correctly identified by their protocol type plus the number of packets m incorrectly identified as not P [17].

$$\text{Recall} = \frac{n \text{ correctly labeled } P}{n \text{ correctly labeled } P + m \text{ incorrectly labeled not } P}$$

Precision is defined as the number of packets n correctly identified as protocol type P divided by the sum of the number of packets n correctly identified as protocol type P plus the number of packets m incorrectly identified as P [17].

$$\text{Precision} = \frac{n \text{ correctly labeled } P}{n \text{ correctly labeled } P + m \text{ incorrectly labeled } P}$$

Finally, we also provide the F1 score for ours and other solutions.

$$\text{F1} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

In addition to these data science measurements, we record, analyze, and compare heap usage using `valgrind`. We also compare the heap allocation for our signatures

and signatures from comparable works to show our solution’s improvements.

We use packet capture (PCAP) data from live network captures for both training and testing scenarios. An interesting note on training data is that in our cases, REXACTOR did not require more than a few hundred starting lines in order to create quality signatures. We found that as long as the starting lines were from a diverse set of captures so that variable URLs, ports, and addresses were not over-represented, the common features were both discovered and generalized enough for re-usability. This demonstrates that REXACTOR may be able to create specific signatures for protocols which may not appear frequently in a data set.

To create HTTP signatures with REXACTOR, we used a small data set of mixed background traffic containing 213 HTTP request packets and 1000 HTTP responses. The resulting signatures and tokens are provided in Table 3.3. For the VoIP protocol tests, we used 252 SIP request packets and 230 responses for training signatures, and 638 RTSP responses and 562 RTSP requests. Our signatures and tokens for these are provided in Table 3.6. We ran tests on data sets from a live capture environment of background web traffic containing 1000 HTTP packets, 1000 SIP packets, 1000 RTSP packets, and 1000 packets of other mixed protocol types, including FTP, SMTP, ICMP, and others.

REXACTOR is a Python tool library designed with Python 3.8. We use the efficient-apriori library (v. 1.1.1) [67] for tokenization, PyShark for data pre-processing (v. 0.4.3) [68], and numpy (v. 1.21) [69] data structures for optimization. The scanning framework used for testing regular expression matching is written in C++ using Boost 1.63 [70] and Hyperscan 5.4.0 [65] libraries for regular expression creation and matching, respectively. All experiments were performed on a single i9 Dell computer with Ubuntu 20.04 installed.

Table 3.3: HTTP Signatures and Tokens

Source	Side	Signatures	Tokens
AutoSig [59]	N/A	N/A	[HTTP/1.]; [GET\0x20/] [HTTP/1.];
LASER [58]	N/A	N/A	[HTTP/1.1]
Wang et al [56]	Request	^(GET HEAD POST).*HTTP/1.(1 0) \x0D\x0A	N/A
	Response	^HTTP/1.(1 0) [2 3 4 5]0.*\x0D\x0A	N/A
LCS Al- gorithm [60]	N/A	^(HTTP/1.1 Host:.u Connection:keep-alive User-Agent:Mozilla/5.0 Accept: Accept- Encoding:ga Accept-Language:en- US.*en;q=0.8 Accept- Charset:ISO-8859-1)	[HTTP/1.1]; [Host:.u]; [Connection:keep-alive]; [User-Agent:Mozilla/5.0]; [Accept:]; [Accept- Encoding:ga]; [Accept-Language:en- US.*en;q=0.8]; [Accept- Charset:ISO-8859-1]
REXACTOR	Request	^(GET / POST / POST POST /mail HEAD /).{0,2048}(HTTP/1.1\r\n)\$	[HTTP/1.]
	Response	^(HTTP/1.1 HTTP/1.).{0,3}.{0,24} \\S(\r\n)\$	[HTTP/1.]

3.4.1 Case Study A: HTTP Detection

The results of scanning related works’ regular expressions and our own on the testing dataset showed comparable recall for Wang et al’s system and REXACTOR. The longest common substring algorithm (LCS) signature performed significantly worse. For our dataset, all systems for both regular expression scanning and literal token scanning had perfect precision, including REXACTOR. Because precision in our data set was ideal for all systems, the F1 Score results are comparable to the recall results.

While precision and recall results were on par with the state-of-the art, REXAC-

TOR’s signature required significantly less heap space than either Wang et al’s solution or LCS when performing regular expression scanning. This demonstrates REXACTOR’s superior performance in real embedded systems where memory optimization is a critical factor.

Table 3.4: HTTP Case Study Results for Regular Expressions

Solution	P	R	F1	Allocation
Wang et al	1.000	0.984	0.992	20,783,011 bytes
LCS	1.000	0.254	0.405	36,015,510 bytes
REXACTOR	1.000	0.976	0.989	18,310,656 bytes

Table 3.5: HTTP Case Study Results for Tokens

Solution	P	R	F1	Allocation
AutoSig	1.000	0.989	0.994	548,715 bytes
LASER	1.000	0.865	0.928	457,778 bytes
REXACTOR	1.000	0.989	0.994	457,280 bytes

3.4.2 Case Study B: VoIP Detection

In addition to a comparison study using HTTP, we ran precision, recall, and memory allocation tests for REXACTOR signatures for two VoIP protocols. In future work, we could expand this to chat, email, and other text-based protocols with the goal of targeting signatures for application-specific signature generation and protocol analysis. For example, we would want to target identifying and distinguishing between Skype, Zoom, or Webex traffic flows [17]. REXACTOR performed amicably with perfect recall and precision in our dataset. Heap allocation results were also comparable to HTTP, showing that the good performance from the case study was not anomalous but rather indicative of good signature construction by REXACTOR.

Table 3.6: VoIP Signatures and Tokens

Protocol	Side	Signatures	Tokens
SIP	Request	^(INVITE sip: ACK sip: BYE sip: REGISTER sip:).{0,74}(SIP/2.0)\$	[sip:];[SIP/2.0]
	Response	^(SIP/2.0 200 OK SIP/2.0 40 SIP/2.0 100 Trying SIP/2.0 180 Ringing SIP/2.0 4 SIP/2.0 18 SIP/2.0) .{0,31}	[SIP/2.0]
RTSP	Request	^(TEARDOWN rtsp:// DESCRIBE rtsp:// SETUP rtsp:// PLAY rtsp://).{0,85}(RTSP/1.0\r\n)\$	[rtsp://];[RTSP/1.0\r\n]
	Response	^(RTSP/1.0)\\d\\d\\d . {0,6}[a-zA-Z](\r\n)\$	[RTSP/1.0]

Table 3.7: VoIP Results for Regular Expressions

Protocol	P	R	F1	Allocation
SIP	1.000	1.000	1.000	34,449,057 bytes
RTSP	1.000	1.000	1.000	14,972,715 bytes

Table 3.8: VoIP Results for Tokens

Protocol	P	R	F1	Allocation
SIP	1.000	1.000	1.000	3,234,815 bytes
RTSP	1.000	1.000	1.000	3,624,176 bytes

3.5 Discussion and Limitations

REXACTOR provides a classification solution for per-packet analysis and does not require data to be read as or reconstructed into flows in order to achieve classification.

The modular design of REXACTOR also allows for components to be used as tool sets so that engineers may adapt the knowledge discovery elements to their own scanning systems or use the full system as a whole framework. For example, the module TRex which derives tokens may be run individually in order to create keywords for Snort filters. Additionally, regular expression signatures can be created and published outside of the scanner in order to be plugged into any other DPI application. REXACTOR also supports tokenization and alignment for raw data, allowing for analysis and signature creation for unknown protocols. This high adaptability means REXACTOR can be applied in many network environments for knowledge discovery and DPI usage. We have generated effective signatures using REXACTOR for email (SMTP, POP3, IMAP), file transfer (FTP), chat (XMP, IRC), VoIP (SIP, RTSP), and web (HTTP) protocol analysis. Because REXACTOR utilizes unsupervised learning and derives features from training data, it could be applied to learning more application specific signatures. One improvement for our solution and others discussed which use unsupervised learning would be an attempt to learn an optimal threshold for feature support rather than rely on manual configuration. In the comparison study, we found that packets not identified by the signatures were HTTP continuation packets of just data payload. Sessionization of protocol flows would allow a detection system to pick up on the first HTTP packet of an exchange with identifying method information, then record the 5-tuple flow information for further scanning. But this technique is not usable in per-packet analysis, and is a limitation of the problem space rather than our solution.

Regular expressions as a text-based approach also have limitations to their applicability in the real, complex network environment. The current industry standard for application layer DPI is to use regular expressions to match expected values in payloads. Protocols or applications may have multiple signatures to match request, response, or message types. Some protocols also do not have strong signatures. Fur-

thermore, different versions of protocols or updated standards and RFCs often require separate signature sets and the continual update of older patterns [71]. Some protocols are also so non-standard that it is difficult to track them on a regular expression signature alone; for example, the payload of an RTP packet is raw data and assigned ports are dynamic [72]. Current state-of-the-art scanners suffer from variable and data-dependent performance impact [73]. With the increasing size and complexity of rulesets and the growing scale of the network environment, regular expression scanning using the current systems is more and more impractical. Regex methods are also limited by what classes they can be used for in network traffic. For example, there is no real regular expression signature which can identify Tor-routed traffic from non-Tor-routed traffic; the data is too diverse. Other desirable classes for traffic, such as traffic type, application, and user or device type present similar issues to regular expression scanners. Each possibility must have its own, human-engineered regular expression to match precisely which is both slow and prone to error. In addition to scalability and machine performance, text-based regular expression matching no longer meets the needs of the environment as over 90% of browsing data today is encrypted [74].

Regular expression matching which relies on deterministic and non-deterministic finite automata (DFA and NFA) is prone to the state explosion problem, where computational complexity increases exponentially as regular expressions increase in number and become more complex. DFAs are fast to search, but are most prone to state explosion as every possibility is explicitly constructed. NFAs provide linear storage space to the size of the regular expression ruleset, but it is relatively easy to arrive at worst-case search performance with rule complexity [75]. Hyperscan is one such tool which constructs an NFA/DFA hybrid for searching. Table 3.9 shows the processing complexity and storage costs of these structures where m is a number of regular expressions of length n .

Table 3.9: Worst-Case Space and Time Complexities for NFA and DFA [1]

Data Structure	Complexity	Storage Cost
NFA	$O(n^2m)$	$O(nm)$
DFA	$O(1)$	$O(\Sigma^{nm})$

Recently, hybrid NFA/DFA engines have been proposed to improve search performance; however, improvements on computational complexity can still be prone to memory issues in the case of large regular expression rulesets [65]. Parallel computing may alleviate some of the computational stress, but is acknowledged as a brute force approach [76]. Hyperscan is able to handle compiling and searching with incredible efficiency; however, this becomes exponentially worse to the point of impracticality when put at a test to scale, which we demonstrate in the next chapter of this work.

CHAPTER 4: ALPINE/PALM - SIMILARITY SEARCH THROUGH LOCALITY-SENSITIVE HASH FINGERPRINTING

4.1 Introduction

The ineffectiveness of existing regular expression-based deep packet inspection methods against encrypted payloads and weakly signed data classes as well as the computational and storage complexity of automata has evidenced the need for other forms of traffic fingerprinting for identification. We propose using locality-sensitive hashing in order to measure similarity of features between packets. Locality-sensitive hashing uses a distance metric to preserve the similarity of original inputs when they are converted to locality-sensitive hashes (LSH). The benefit of locality-sensitive hashing here is that it is capable of reducing high-dimensionality data into a low-dimensional, space-effective, computationally efficient representation of the same data while still preserving similarity metrics between the data points [77]. Hash families may be defined for many distance measures such as Hamming distance, L1 and L2 norm [78], and Jaccard similarity [79].

While cryptographic hashes attempt to minimize collisions, locality-sensitive hashes preserve the similarity of data points by generating hashes closer to one another depending on the input's similarity. The following definition intuitively states that data points which are locally nearby have a higher probability of collision than further points.

Definition 1 [2] *A family H of functions from a domain S to a range U is called (r, ϵ, p_1, p_2) -sensitive, with $r, \epsilon > 0$, $p_1 > 0$, $p_2 > 0$, if for any $p, q \in S$, the following conditions hold:*

- if $D(p, q) \leq r$ then $Pr_H[h(p) = h(q)] \geq p_1$,
- if $D(p, q) > r(1 + \epsilon)$ then $Pr_H[h(p) = h(q)] \leq p_2$.

ALPINE [75] is **A** **L**ocality-Sensitive **P**acket **I**nspection **E**ngine which performs shallow packet inspection on selected features from the network and transport layer encapsulation headers of packets. PALM performs **P**ayload **A**nalysis using **L**ocality-Sensitive **M**easurements. Both these models use locality-sensitive hashing as a means for normalizing data into an embedding which can then be compared for similarity and classified using distance metrics.

Locality-sensitive hashing has been used in user-level browser fingerprinting, where web-based browser fingerprints are created by extracting multiple values from the browser API and hashed using MinHash or some similar algorithm. This generates signatures of high entropy, where the data is uniquely identifiable as a particular device or host [80]. LSH clustering has also been used to make signatures for classifying malware at scale [81]. In a similar fashion, we apply this technique to network packet data by using features which should evidence values unique to a given class and creating the locality-sensitive hash from those.

Similarity search approximation can be performed by linearly comparing the hashes created from the packets. Similarity search is a generalized term in data mining which refers to searching for objects where the available comparator is some common pattern or set of similarities. This group of mechanisms includes tools such as Nearest Neighbors searching, link-based similarity searches, duplicate detection, and object representation [2]. This can be applied in cyber criminal investigation to find individuals in the same criminal network, content similar to known evidence, images similar to another, or known fingerprints across network data. Qualities of good similarity search methods include efficient querying and storage mechanisms, a minimized memory footprint, and minimal human interference [82].

One data engineering challenge is constructing suitable, unique *tokens* which characterize the data meaningfully. The sentence, “Palm trees are native to the Pacific”, may be split on whitespace into the set of tokens $T = \{\text{“palm”}, \text{“trees”}, \text{“are”}, \text{“native”}, \text{“to”}, \text{“the”}, \text{“pacific”}\}$. It is obvious that tokens like “the” are far too generic in the English language to be characteristic of this sentence, but a token like “palm” or “pacific” will be much more unique and indicate a stronger similarity. Formally, this idea is quantified as taking the inverse document frequency (IDF), where we minimize the importance of terms which appear frequently in the document set, and combining it with term frequency (TF) as the TF-IDF value [83]. This value describes tokens which are both common and characteristically unique of the document set.

4.2 Previous Work

Jaccard similarity measures the intersection of two sets over their union. This value is often used in recommender systems, and can also be used to detect plagiarism or make predictions. In order to take the Jaccard similarity, data must be split up into sets, which are made up of individual elements like the tokens or shingles of strings as done in SigBox [45] and REXACTOR [71] for regular expression signature generation.

Some limited work has been done in application fingerprinting through locality-sensitive hashing. Tang et al proposed HSLF [84], an HTTP header sequence-based LSH fingerprint generator for classifying applications in HTTP traffic. Their results show the ability of their SimHash-based method to distinguish accurately among data such as VMWare, Firefox, and WeChat. This work is limited as it only applies to cleartext HTTP traffic. Furthermore, the scope of this traffic only makes up a fraction of real-world data.

Jiang and Gokhale [85] use locality-sensitive hashing on packet-level features to accurately classify network traffic statelessly. From a model perspective, their use of K-Nearest Neighbors clustering scales sub-optimally. Their work also focused solely on multimedia applications versus legacy web-browsing; there is no expansion into traffic

type, application, protocol, or other needed classification in current DPI. ALPINE and PALM expand these works into a much more complex and diverse set of experiments and applications.

4.3 Methodology

4.3.1 MinHash and Locality-Sensitive Hash Forest

In order to form the elements of a set, we *shingle* packets into items and reduce the data to a set intersection problem. An object is w -shingled when a sub-sequence of length w of contiguous tokens is cut from it. In ALPINE, we use the port information, transport layer protocol, flag values, and packet length to make up our shingles. In PALM, word tokens are created by delimiting packet payloads by whitespace. Thus, each original packet P is now associated with some set S_P of shingles as depicted in Figure 4.1.

$$J(P_A, P_B) = \frac{|S_A \cap S_B|}{|S_A \cup S_B|} \quad (4.1)$$

Figure 4.1: Jaccard Similarity of Shingle Sets

Jaccard similarity can be approximated quickly through the MinHash algorithm. The hash representation also fixates and potentially reduces the amount of space packet data occupies, and furthermore normalizes the data in cases where the shingles are heterogenous and of varying length and type. The MinHash of a given column is the number of the first row in permuted order whose characteristic matrix value is 1. The sequence is continued down the columns and repeated for k permutations. The probability that the MinHash function for a random permutation of rows produces the same value for two sets is a close approximation to their Jaccard similarity [86]. Bawa et al [2] developed the LSHForest algorithm which indexes the locality-sensitive hashes in a forest-based data structure in order to optimize the similarity search process. It improves upon nearest-neighbor searches by both reducing complexity and eliminating

the need to know the distance r from the query to the nearest neighbor of a given data point. Optimization is achieved through the use of a specific set of hash functions which ensures that any nearest neighbor query returns ϵ -approximate neighbors so long as a suitable l number of trees is chosen. LSHForest scales much more efficiently compared to previous work using nearest neighbor searches [85]. The build time of a KNN model can have a complexity of $O(n^2)$, where n is the number of items. In contrast, LSH construction is done in linear time [2]. The theoretical complexities of operations on LSHForest are provided from the original paper in Table 4.1. When compared with time and space complexities for the automata in regular expression searching in Table 3.9, the logarithmic improvement and fixed storage space is clear. We chose to use this data structure in the classification stage of ALPINE and PALM in order to utilize this optimization.

Table 4.1: Time Complexities for LSHForest [2]

Operation	Complexity
Insertion	$O(l * \log_B n)$
Deletion	$O(l * \log_B n)$
Query	$O(l * \log_B n) + O(l * \log B) + O(M/B)$

In these equations, l represents the number of trees, n the number of data points in the dataset, and B the branching factor of an internal node. Storage is optimized to be linear, $O(n)$, through the use of compressed PATRICIA tries [2]. The m nearest neighbors for some point q is found by first performing a binary, longest matching prefix search for a leaf node with the point s_i . Then some M points are collected synchronously across all trees and ranked by similarity score, with the top m being returned as an answer. The m number is used as votes to classify the sample by majority once the query is complete. This majority vote approach also allows for multiple classification results; for example, an adaptive system may be able to take the second closest result if the first choice turns out to be incorrect.

4.3.2 ALPINE: A Locality-Based Packet Inspection Engine

ALPINE uses Jaccard similarity in order to determine likeness between sets of the following packet header features: *source and destination IP address, type of service field, packet length, IP next protocol, IP flags, source and destination port* (see Figure 4.2). For indexing the computed MinHashes and performing lookup for classification, ALPINE uses the previously discussed LSHForest [2]. A query is performed by computing the MinHash of the input data and doing a binary search on the prefix trees. The hash is then mapped back to its label index. A majority vote is taken from the returned m nearest neighbors and the data is classified accordingly.

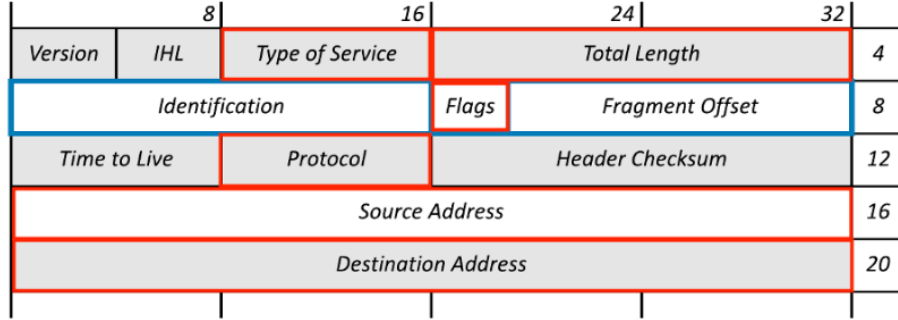


Figure 4.2: Example TCP/IPv4 header with extracted features for ALPINE outlined

4.3.3 PALM: Payload Analysis Using Locality Sensitive Measurements

PALM employs a similar strategy to previous work in payload signature generation like SigBox [45] and RExACtor [71]. PALM delimits packet payloads by whitespace in order to create word-style tokens. In future work, other delimiters or a sliding-window shingling technique may suit identification of certain data types such as binary application layer protocols like HTTP2. In the experiments presented in this report, whitespace proved to be the most effective delimiter. MinHash was used to generate the locality-sensitive hashes for the data samples and test samples. We use the LSHForest algorithm developed by Bawa et al [2] to index the locality-sensitive hashes and optimize the search process.

4.4 Experiments and Results

Table 4.2: Description of Combined Dataset

Application Type	Protocol	# Samples
<i>File Transfer</i>	FTP	10,015
	FTP-DATA	4,000
	BitTorrent	20,648
<i>Voice-over-IP</i>	MGCP	1,568
	SIP	1,112
	H.225	1,300
	RTP	15,552
	RTCP	1,626
<i>Mail</i>	SMTP	5,981
	POP	1,675
	IMAP	3,318
<i>Authentication & Access</i>	LDAP	1,354
	SMB	3,554
	Telnet	1,888
<i>Tunneling</i>	GPRS	9,981
	PPTP	1,288
	SSH	3,039
<i>Web & Chat</i>	DNS	10,563
	HTTP	21,298
	XMPP	1,553
	TLS	4000
<i>Networking</i>	DHCP	1,444
	NBNS	1,216
	GQUIC	1,740
	NTP	1,940
	SSDP	8,504

Using only a single dataset or network environment can lead to bias in results [87]. Thus, we created a combined dataset from many available public datasets in order to diversify the traffic and contents. Our sources include the CDX 2009 captures [88], the Skynet Tor dataset [89], the Canadian Institute for Cybersecurity’s ISCX VPN/non-VPN and Tor/non-Tor datasets [90, 91], and repositories from Wireshark, Cloudshark, and IEEE Dataport. This repository contains PCAPs from twenty-six different, common application layer protocols. For experimental reproducibility and general use we

have made this new dataset available publicly ¹.

4.4.1 Protocol Auto Detection

It is crucial for lawful interception and DPI middleware boxes to be able to autodetect application-layer protocols in order to properly assess, parse, and process seen traffic. The environment on any given signal can be wildly diverse; thus, classification systems must be capable of scaling to large, multiclass problems. Previous systems [84, 85] have not attempted protocol detection to the scale of a 26-class problem, making this model a true and fresh pioneer in the DPI field. We first split our data into training and test sets and extracted features from each of the protocols. Next, we generated hash values for the training samples and indexed them into the forest. For testing, we generate a hash of the input value and perform the similarity search on the trained model. The results of the testing samples are provided in Table 4.3. The evidence shows that the combined model has high accuracy and precision/recall for protocol identification. Cohen’s kappa also shows a strong agreement among the ensemble voters. The method also performs equally well across application types, indicating generalizability. Figure 4.3 shows the confusion matrix results of the classifier’s ensemble voting system.

4.4.2 Traffic Type

Standardized documentation necessitates that there are some detectable similarities between protocols. It is reasonable to think that detection of more abstract classes like traffic type (i.e. email, file transfer, web data) would be more heterogeneous and thus more difficult to detect by embedding and similarity search. We experimented with classification by traffic type by amalgamating the dataset into classes as described in Table 4.2 in the datasets section of this work. Results in Figure 4.4 show a confusion matrix of the classification results which are also highly accurate.

¹<https://github.com/mayakapoor/protocol-dataset>

Table 4.3: Protocol Autodetection Results

	Jaccard			TF-IDF		
Accuracy	0.996			0.841		
Cohen's κ	0.997			0.841		
Protocol	P	R	F1	P	R	F1
BitTorrent	1.00	0.99	0.99	0.89	0.85	0.86
DHCP	1.00	1.00	1.00	0.91	1.00	0.95
DNS	1.00	0.99	1.00	0.95	0.97	0.96
FTP	0.99	1.00	0.99	0.83	0.90	0.86
FTP-DATA	1.00	1.00	1.00	0.93	0.96	0.94
GPRS	1.00	1.00	1.00	0.91	0.98	0.95
GQUIC	1.00	0.99	0.99	0.87	0.75	0.81
HTTP	1.00	0.99	0.99	0.82	0.48	0.61
H.225	1.00	1.00	1.00	0.83	0.97	0.90
IMAP	1.00	0.99	0.99	0.86	0.79	0.82
LDAP	1.00	0.99	0.99	0.87	1.00	0.93
MGCP	1.00	1.00	1.00	0.90	0.65	0.75
NBNS	0.99	1.00	0.99	0.94	1.00	0.97
NTP	1.00	1.00	1.00	0.91	1.00	0.95
POP	0.98	1.00	0.99	0.91	0.47	0.62
PPTP	0.95	1.00	0.98	0.84	0.90	0.87
RTCP	1.00	1.00	1.00	0.96	1.00	0.98
RTP	1.00	1.00	1.00	0.89	0.88	0.89
SIP	1.00	1.00	1.00	0.54	0.98	0.69
SMB	1.00	0.97	0.98	0.95	0.99	0.97
SMTP	1.00	0.99	0.99	0.92	0.93	0.93
SSDP	1.00	1.00	1.00	0.72	0.28	0.40
SSH	1.00	0.98	0.99	0.90	0.86	0.88
Telnet	0.99	1.00	0.99	0.86	0.99	0.92
TLS	1.00	1.00	1.00	1.00	1.00	1.00
XMPP	1.00	1.00	1.00	0.79	0.90	0.84

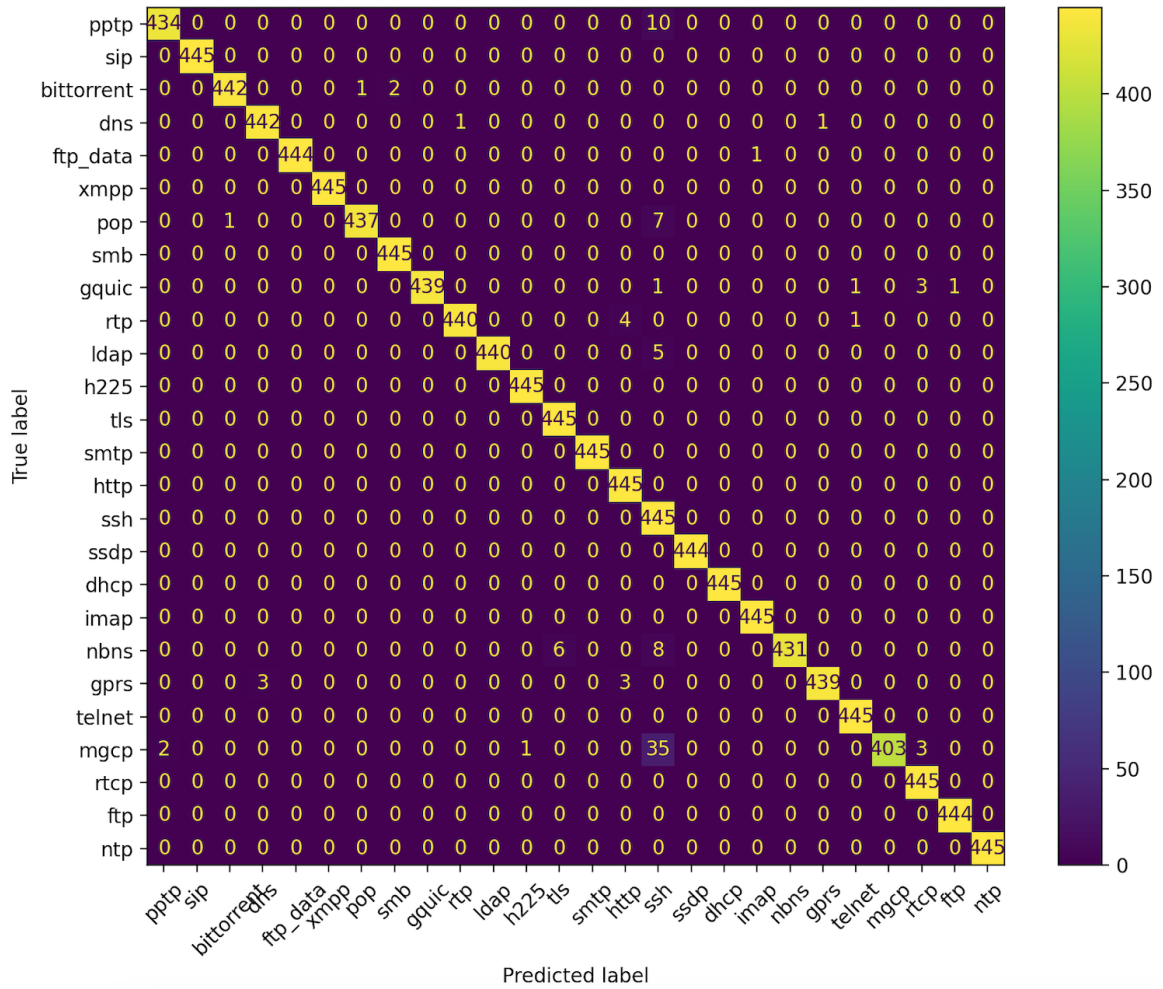


Figure 4.3: Twenty-six Class Identification Problem using the ALPINE-PALM Combined Model

4.4.3 TF-IDF Score Evaluation

There are many complex ways to measure similarity between two items; Jaccard similarity is considered relatively simple. In order to reason whether or not the similarity metric was effective, we used the TF-IDF score of payload tokens in order to filter out low-value tokens and generate hashes only from high-value ones. This metric was chosen to address concern that large data payloads, for example HTML documents or email messages, might introduce noise into hashes which would reduce classification accuracy. Instead, isolating payload tokens to only the high value ones could reduce spurious extra data. The TF-IDF measurement of tokens was imple-

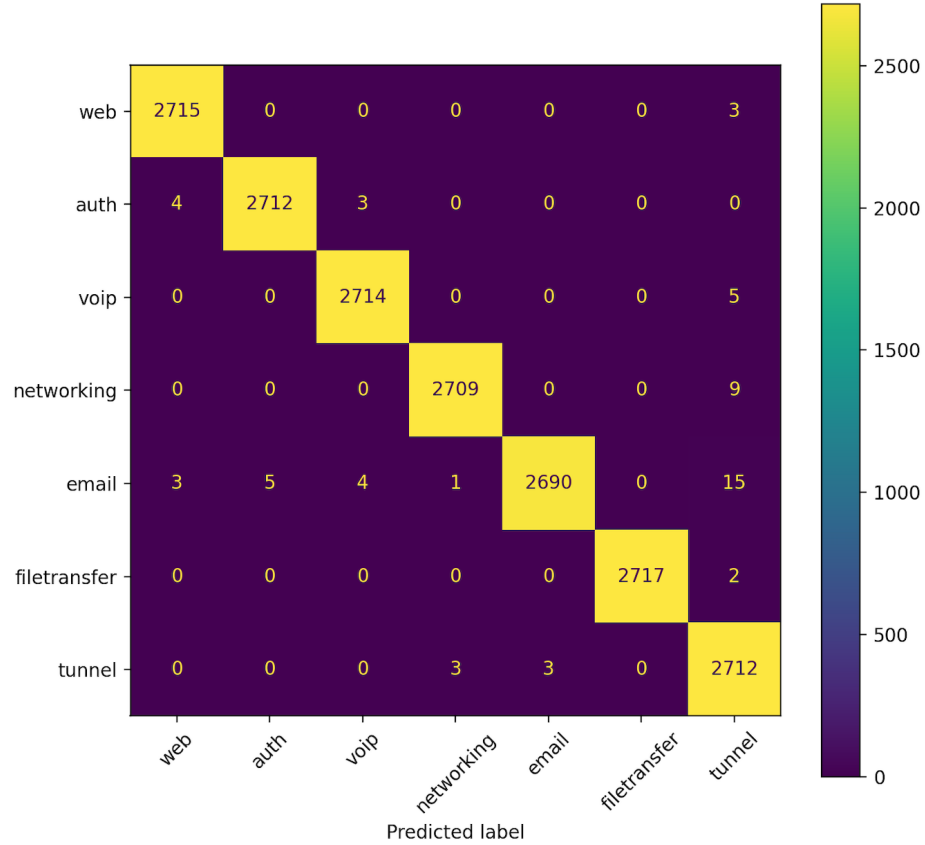


Figure 4.4: Confusion Matrix of Traffic Type Classification Results

mented in the training stage by modifying configurations of the PALM model to filter the added tokens. Only those tokens with TF-IDF score above the minimum threshold were kept and included in the actual LSH signature. The results in Table 4.3 show that this actually had a negative impact on the classification accuracy, indicating the basic Jaccard similarity was a better distance metric for this use case.

4.4.4 Multi-Embeddings

It has been shown that concatenating multiple classifiers and evaluating the agreeance among voters yields a reduction in error and even correlation when there is a misclassification [92]. We hypothesized that enabling both ALPINE and PALM to run together would yield the most accurate and agreed-upon classification results. We tested this theorem by running the 26-class protocol identification test against con-

Table 4.4: Multi-Hash Embedding Classification Results

	ALPINE			PALM			Multi-Model		
Accuracy	0.988			0.717			0.996		
Cohen's κ	0.988			0.705			0.997		
Protocol	F1	P	R	F1	P	R	F1	P	R
BitTorrent	0.98	0.97	1.00	0.74	0.75	0.73	0.99	1.00	0.99
DHCP	1.00	1.00	1.00	0.91	0.83	1.00	1.00	1.00	1.00
DNS	1.00	1.00	1.00	0.24	0.46	0.17	1.00	1.00	0.99
FTP	0.96	0.97	0.94	0.90	0.83	0.98	0.99	0.99	1.00
FTPDATA	1.00	1.00	1.00	0.20	0.37	0.13	1.00	1.00	1.00
GPRS	1.00	1.00	1.00	0.91	0.98	0.85	1.00	1.00	1.00
GQUIC	0.99	0.99	1.00	0.27	0.48	0.19	0.99	1.00	0.98
HTTP	0.99	1.00	1.00	0.61	0.65	0.55	1.00	1.00	0.99
H.225	1.00	1.00	1.00	0.63	0.68	0.95	1.00	1.00	0.99
IMAP	0.99	0.94	0.98	0.85	0.77	0.59	0.99	1.00	0.99
LDAP	0.96	1.00	0.98	0.66	0.75	0.59	0.99	1.00	0.99
MGCP	1.00	1.00	1.00	0.90	0.83	1.00	1.00	1.00	1.00
NBNS	0.99	0.98	1.00	0.87	0.82	0.92	0.99	0.99	1.00
NTP	1.00	1.00	1.00	0.90	0.82	0.99	1.00	1.00	1.00
POP	0.99	0.98	0.99	0.26	0.41	0.20	0.99	0.98	1.00
PPTP	0.99	0.98	0.99	0.49	0.33	1.00	0.98	0.95	1.00
RTCP	1.00	1.00	1.00	0.92	0.85	1.00	1.00	1.00	1.00
RTP	0.99	0.98	0.99	0.40	0.65	0.29	1.00	1.00	1.00
SIP	1.00	1.00	1.00	0.93	0.87	1.00	1.00	1.00	1.00
SMB	0.99	1.00	0.98	0.81	0.80	0.82	0.98	1.00	0.97
SMTP	0.98	0.99	0.98	0.81	0.80	0.82	0.99	1.00	0.99
SSDP	0.99	1.00	0.98	0.92	0.86	1.00	1.00	1.00	1.00
SSH	0.95	0.97	0.92	0.47	0.67	0.37	0.99	1.00	0.98
Telnet	0.99	0.98	1.00	0.80	0.75	0.85	0.99	0.99	1.00
TLS	0.99	0.99	1.00	0.61	0.51	0.75	1.00	1.00	1.00
XMPP	1.00	1.00	0.99	0.92	1.00	0.85	1.00	1.00	1.00

figurations of the model with just ALPINE embedding, only PALM embedding, and a final combined run. Results as shown in Table 4.4 demonstrate that the combined votes of both models were the most accurate in the final tests. Furthermore, the agreeance among voters was also highest using multiple hash embeddings.

4.4.5 Model Performance and Throughput

Applied machine learning must keep up with signal processing at very high rates; furthermore, any passive system must not interfere with legitimate traffic and service and can thus not introduce additional latency or the potential for dropped packets. The software must also be capable of performing the required processing for data analysis on as much of the traffic as possible (ideally, all of it). Diverting or copying traffic offline for analysis as well as parallelization strategies and traffic thinning and load balancing are legitimate system design choices which can mitigate this. But there is still room for improvement and optimization at the classification solution level to achieve real-time rates. Thus, we perform experiments to see how the system scales based on model sizes and number of classes. We measured the training time, system throughput in millibits per second (Mbps) as well as memory usage during the testing phase in megabytes (MiB) as we were running all performance tests on the combined dataset consisting of 140,157 total packets. Random under-sampling was used to even out the distribution of data and avoid bias, which caused our total number of packets after sampling to be a fraction from the original set. For testing the number of classes, we created experiments with binary RTP/non-RTP classification, the traffic-type multiclass experiment with 8 classes/types, and the largest 26-class experiment where we identify all possible data layer protocols. In Table 4.6, we detail the results of the experiments for a varied number of classes and sample sizes.

In order to test system scalability in deployment and verify the theoretical complexity discussed in the background section of this work, we performed a series of experiments increasing the number of signatures for a model used to classify HTTP

Table 4.5: Performance Results for Varied Number of Classes

Classes	# Hashes	# Samples	Training Time	Test Time	Memory Usage	Throughput	Accuracy
RTP / non-RTP <i>Binary classification</i>	20,613	13,743	4.273s	1.21ms	470.8 MiB	2.596 Mbs	1.000
	2,400	1,600	0.738s	0.866ms	283.7 MiB	3.139 Mbs	0.997
	240	160	0.075s	1.093ms	353.4 MiB	2.317 Mbs	1.00
Traffic Type <i>8-class problem</i>	28,543	19,029	7.136s	1.015ms	584.6 MiB	3.728 Mbs	0.997
	12,600	8,400	3.31s	1.006ms	305.3 MiB	2.842 Mbs	0.997
	240	160	0.683s	2.397ms	290.6 MiB	2.54 Mbs	0.998
Protocol <i>26-class problem</i>	17,347	11,565	5.312s	1.138ms	352.4 MiB	3.041 Mbs	0.997
	15,600	10,400	3.13s	0.727ms	305.3 MiB	4.231 Mbs	0.998
	1,560	1,040	0.492s	0.901ms	301.8 MiB	3.212 Mbs	0.994

traffic. For a comparative baseline, we used the regular expression scanning tool built in REXACTOR [71] which uses the Hyperscan 5.4.0 regular expression matching library by Intel [65]. All experiments were performed on a single i9 Dell computer with Ubuntu 18.04 installed. In the trials, the number of signatures represented the number of LSH hashes our system or the number of regular expressions added to the sniffer application. We used REXACTOR [71] to generate HTTP signatures from the HTTP traffic in our dataset. Both models performed the classification with 100% accuracy on all trials.

4.5 Discussion and Limitations

The locality-sensitive hashing approach is an improvement upon the regular expression approach to DPI in several ways. First, the hash generation process reduces

Table 4.6: Performance Results for ALPINE and PALM as a Multi-embedding Versus Hyperscan

System	# Signatures	Training	Testing	Memory
ALPINE/PALM	1,000	1.21s	1.23s	392.4 MiB
	10,000	2.75s	3.56s	419.9 MiB
	100,000	52.75s	4.2s	1047.4 MiB
	1,000,000	423.3s	4.52s	4717.9 MiB
Hyperscan	1,000	4.6s	156ms	70.7 MiB
	10,000	81.8s	1.3s	529.3 MiB
	100,000	39.9m	17.28s	4596.7 MiB
	1,000,000	153m	63.1s	-

the need for expert knowledge about a particular protocol. Previously, engineers would have to employ unique signatures for different message types. It would not be uncommon, for example, to have regular expression signatures for every method for a given text-based protocol like FTP or IMAP. Furthermore, requests and response messages often require unique signatures. Regexes are difficult for humans to read and comprehend. The management of server-client side interactions is also a challenge to know which versions to apply. Instead, hash embeddings using our forward-thinking system create a unique value for each provided input and indexes those values into one bucket. This even more generally captures the broad diversity of possible payload and header contents while providing a new layer of abstraction and simplicity. Flows may be examined bi-directionally, and all message types and methods considered. We are also not limited to traffic flows and protocol types. As demonstrated by the experiments in this section, the system is capable of distinguishing more heterogeneous or abstract class types like traffic type, mechanisms like encryption, and traditionally weak-signed traffic like RTP. Furthermore, a single bad bit can cause a packet match on regular expression to fail. Network traffic and packets in the wild are extremely diverse and can be prone to error, and thus are an ideal test set for the multi-embedding hash technique which is both more informative and resilient. In terms of performance, the LSHForest as implemented in our model further reduces

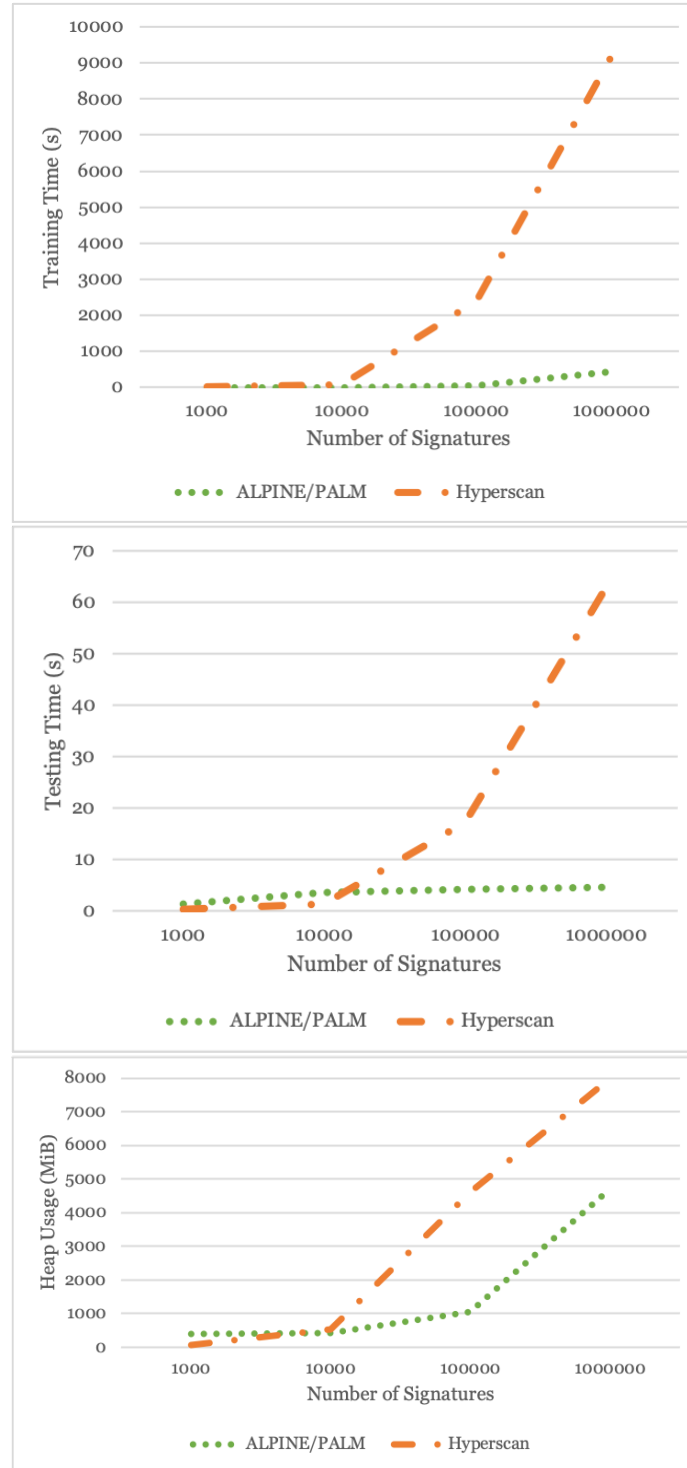


Figure 4.5: Time and Space Complexity Results for ALPINE/PALM versus Hyperscan

the storage and computational complexity when compared with a regular expression, automata-searching approach.

We posed the following research questions to assess the ALPINE and PALM models and report our findings here:

- **R1.** Can packets be classified by application-layer potocol through the creation and index of locality-sensitive hashes?

Our model performed at a 99.6% accuracy rate in its identification of many application layer protocols, indicating that locality sensitive hashes may be constructed which are highly representative of the data and low-cost for searching when indexed. This model may be expanded in real world implementations of DPI for traffic surveillance, cybersecurity and intrusion detection, quality of service managment and network management, content and copyright enforcement, and traffic profiling. Detecting the application layer protocol remains an important task for engineers. This hashing technique improves upon the available state of the art in both accuracy and breadth of scope in application.

- **R2.** Does this model generalize to other kinds of network traffic classes besides protocol? Could we also distinguish more variable data such as email or VoIP traffic classes?

In addition to the application layer protocol, we also used traffic type labels on the dataset for experimentation. In this problem also the hash embedding technique proved highly accurate, which further demonstrates the generalizability of the models to new problems and real data.

- **R3.** Does a different word embedding technique such as using only high-value tokens as characterized by TF-IDF score improve distance measure results? Will this reduce noise in the data and isolate important features?

We used Jaccard similarity as a baseline for comparison, but found that this distance metric was more effective than the TF-IDF score. This highlights a

counterintuitive but significant point that challenges current network analysis trends. Machine learning techniques have become more technically advanced and often computationally intensive in order to improve model accuracy or other desired metrics. While TF-IDF decently diminishes tokens such as “a” and “the” from the English language, the fact that “HTTP\1.1” appears in nearly every HTTP request in our dataset is actually relevant considering the context. Especially in network traffic analysis where scalability and line-rate capability is highly desirable, simple techniques can still be effective and should not be disregarded. Our experiments evidence that a more complex model is not always the appropriate direction. Even within subfields of machine learning like natural language processing, there is not a silver bullet solution which triumphs in every situation. Rather, having multiple models to try or angles from which to consider the data seems to provide better, more well-rounded results.

- **R4.** If we combined multiple hash embeddings, i.e. both ALPINE and PALM, and concatenate their results, will there be an improvement in classification performance over a single hash embedding of the feature set?

Multi-embeddings are a proposed solution particularly for more heterogeneous datasets. This captures multiple angles or representations of the data and combines them for a more full-scope picture and analysis. In contrast, regular expressions provide only a one-dimensional viewpoint and are incapable of encoding highly variable data without impractical complexity and data over-selection. The modular design of our system enables using one or more embeddings, which allows for a balance between performance impact and accuracy. Interestingly, though the PALM model was overall poorer-performing than the ALPINE model, votes from the PALM model helped improve the combined result regardless.

- **R5.** Does the sample size, number of hashes, or number of classes significantly

impact throughput or accuracy?

The performance impact of applied machine learning is a major concern for traffic analysis and cybersecurity. We provide transparent numbers of system performance to present the solution at scale. All experiments in this trial were run on a 1.6 GHz Dual-Core Intel Core i5 processor with 16GB DDR3 memory. In future work, leveraging an accelerated processor such as an FPGA would greatly enhance performance and throughput due to the repetitive nature of hash computations. The most impactful procedure based on system profiling is the MinHash step, which is likely optimizable through an FPGA [85]. The measure of acceptable throughput depends on the network environment; however, in terms of scalability it is valuable that our model is data independent in that it does not seem to increase exponentially or vary in performance wildly depending on the data. Our model also trains quickly, meaning that models could be swapped or trained online or in the field if needed. Future research could further optimize the process through batching and parallel processing. For example, it would be possible and beneficial to run training computations in parallel. Distance measures could also be taken in parallel for different models on separate threads and then the results aggregated.

- **R6.** What is the rate of scalability of the system? Does it outperform state-of-the-art regular expression scanning in terms of testing/training time and memory usage?

We indicate in Table 4.6 that for small signature sets, Hyperscan initially outperformed our model; however, when scaled to larger signature sets, our model clearly is much more capable of handling the extra load. The comparative line graphs in Figure 4.5 illustrate the exponential versus sublinear growth patterns of each measure of performance, clearly demonstrating our systems' scalability

over the state of the art. In constructing models with 1 million signatures, our model takes just over 7 minutes while the Hyperscan implementation takes over 2 hours. At scale, even our testing time is much smaller than the baseline. Furthermore, at the 1 million signature scale our model uses just over 4GB of heap space, while Hyperscan crashed our setup due to exceeding 8GB of dynamic memory for 1 million signatures. This illustrates the point made earlier that no single model is best suited for every situation; rather, multiple models and strategies are needed for unique network environments. If the problem set or data structure size is small, Hyperscan could be an optimal choice; however, our model is more generalizable to larger data sets and adaptable to a variety of problems. It would not be unreasonable to scale to several million signatures in a real surveillance or DPI system which ALPINE and PALM would be better suited for than the current capability of regular expression matching.

One key takeaway from the research done in ALPINE and PALM is that multi-embeddings on packets in our given datasets and experiments were more effective at classifying the data than any single form of analysis. Naturally, this comes at the cost of generating more embeddings and thus increasing complexity and latency. Thus, when designing a multi-embedding model, there is a balance to be achieved between number of layers of representation and desired results. If storage space of models and embeddings or CPU utilization is of more concern than acute accuracy of classification, fewer or less complex representations may be required. The data may also be aptly divided for input into different models; for example, ALPINE considers header features while PALM considers the payload.

ALPINE and PALM show the potential for locality-sensitive measurements to consider multiple layers of features and are resilient to non-standard implementation. They normalize heterogeneous data into fixed-size, comparable space. These models as part of a DPI system have potential to scale both in performance and storage ca-

capacity. But, there are still problems in network classification that elude this method of feature extraction. For example, encrypted or compressed payloads make the text-based approach of PALM or regular expressions generated by REXACTOR ineffective. Distinguishing encrypted traffic payloads from plaintext or compressed payloads is also a classification problem not appropriately dealt with through these methods. If ALPINE features of the header are not enough alone to distinguish the traffic type or particular application the traffic represents, then the classification will be inaccurate. If the header is encrypted, this also renders ALPINE ineffective. In further research, a method of understanding and profiling encrypted and compressed traffic or obfuscated payloads would greatly enhance deep learning-based DPI systems.

CHAPTER 5: MAPLE - IMAGE REPRESENTATION AND RECOGNITION OF INTERNET TRAFFIC

5.1 Introduction

In chapters 3 and 4, we discuss the difficulty of generating token-based or regular expression-based signatures for weakly signed, variable data. This type of data floods modern Internet traffic today; examples would be Voice-over-IP (VoIP), streaming services, downloads and file transfers, and peer-to-peer sharing. Furthermore, encrypted traffic payloads or compressed data make generating signatures from payload text ineffective. The tokenization strategy of PALM is defeated by these techniques, and as ALPINE focuses only on header features there is still the problem of an unanalyzed payload with potentially useful, untapped information.

In 2021, it was projected that 3 billion people use VoIP technology as a regular part of their daily lives. In 2021-2022, Zoom reported over 300 million users [93], Webex reported over 600 million, and Microsoft Teams recorded 275 million participants [94]. The widespread adaptation of VoIP technology has highlighted several cybersecurity vulnerabilities which make this type of traffic even more important to network intrusion detection systems. It is simple to spoof IP addresses or Session Initiation Protocol (SIP) universal resource identifiers with VoIP technology and produce robocalls used in phishing/spam or in denial of service attacks [95]. VoIP packet data is vulnerable to being the transport for backdoors, worms, trojans, and viruses which can be embedded [96, 97]. Cybersecurity specialists designing intrusion detection systems need to be able to intercept and process VoIP streams in order to spot these kinds of intrusions [98]. VoIP call analysis and reconstruction can also be used in the forensic environment to provide critical evidence of criminal activity. Policy

or content rights violations or copyright infringement may also be detected through the analysis of the reconstructed calls [99, 100]. This process and search capability is also useful to intelligence operations to find mission critical data from the enormous amount of traffic flowing in mission networks [101].

A majority of mainstream VoIP protocols use the real-time transport protocol (RTP) as their transport layer for encapsulation. The application layer protocol, for example the Session Initiation Protocol (SIP), Media Gateway Control Protocol (MGCP), or H.248, will establish the connection between endpoints and carry important session information such as codec encoding and metadata for the call. The application then relies on RTP to manage the dataflow between the connected systems. In the *complex, real network environment*, the implementation of SIP and RTP data flow is often sent de-coupled across physical signals and ports which presents a challenge for cybersecurity specialists using middlebox technologies for packet inspection and call reconstruction. Because the RTP data is encoded, it is necessary to know the signaling information in order to retrieve the codec information for proper decoding. Furthermore, the metadata associated with signaling protocols such as the SIP URI identifier provides necessary enrichment and context for the actual call. Lastly, RTP data may arrive at the middlebox before the signaling information which contains the port numbers associated with the RTP stream for that particular call; thus, the middlebox must retrospectively identify data packets (RTP) which belong to particular signaling information headers (SIP, for example).

Detecting the RTP stream itself from other types of traffic can be a difficult problem because the protocol's signature is weak and can be further obfuscated by encryption (SRTP). The problem of accurately classifying network traffic has expanded beyond the scope of capability of text-based solutions. Instead, we propose using higher dimensionality of network traffic in order to advance RTP detection capability. MAPLE is a **MA**trix-based **PayLoad Encoder** which transforms packets into grayscale images

to unveil hidden representation which may be used as input into a convolutional neural network model. Using this approach, we are able to expand detection capability to weakly signed data, encrypted versus compressed traffic, and better traffic type and application profiling.

5.2 Previous Work

Data in convolutional neural networks are organized into three aspects which determine the spatial dimensionality of the input: the height, width, and depth. The data engineering design processes pixelated data and discovers visual patterns in the latent space.

Following the input layer, there are three layer types which make up the CNN model. A *pooling* layer performs down-sampling along the spatial dimensionality of the given input in order to reduce dimensions. For each feature map $a_{i,j,k}^l$, pooling can be calculated as follows:

$$y_{i,j,k}^l = \text{pool}(a_{m,n,k}^l), \forall (m,n) \in R_{ij} \quad (5.1)$$

where R_{ij} represents the local neighborhood around point (i,j) . Common pooling functions include max pooling and average pooling [102]. A *convolutional layer* uses a kernel filter in order to compute feature maps. The feature value of location (i,j) in the k th feature map of the i th layer $z_{i,j,k}^l$ is calculated by:

$$z_{i,j,k}^l = w_k^{lT} x_{i,j}^l + b_k^l \quad (5.2)$$

where w_k^l and b_k^l are the weight vector and bias term of the k th filter at the l th layer. The activation function, usually sigmoid, tanh, or rectified linear unit (ReLU), may then be calculated as:

$$a_{i,j,k}^l = a(z_{i,j,k}^l) \quad (5.3)$$

Finally, *fully-connected* layer takes input from the previous layer and connects them to neurons of the output layer in order to perform high-level reasoning or learn some global semantic information [103].

LeNet [104] consists of two sets of convolutional, activation, and pooling layers. There is also a final set of consisting of fully-connected, activation, and fully-connected layers which proceed to a softmax normalization and classification at the output layer.

Deep neural networks suffer from accuracy degradation due to how difficult it can become to map intermediary or residual layer outputs to inputs. Thus, a residual convolutional neural network (ResNet) is built up of blocks of stacked layers formed from a series of convolutions and non-linear activation functions. Following the authors who first introduced deep residual networks [105], we define a block as:

$$y = F(x, \{W_i\}) + x \quad (5.4)$$

Where x and y are the input and output vectors, respectively, and $F(x, \{W_i\})$ is the residual mapping with a set of weights W_i . If $H(x)$ represents the ideal output mapping that corresponds to the ground truth, residual networks hypothesize that $F(x) + x = H(x)$. The residual function is pushed to zero such that the equation becomes an identity function $H(x) = x$. This identity mapping, or shortcut layer, is what reduces the degradation problem and minimizes training error as the network becomes deeper [105].

In previous work, CNNs have been employed to classify network traffic both for per-packet and per-flow scenarios. Lim et al [106] also transform packets into grayscale images and use CNN and ResNet architectures to label them by application (RDP, Skype, BitTorrent, and others). When using payloads of 1024B and a ResNet architecture, they achieved a 0.97 F1 score for accurately classifying the packets as one of eight types of applications.

Yang et al [107] transform network packets in creative adversarial networks (CANs) using quantile normalization which generate pixelated images for CNN processing. They combine multiple packets into a single image for flow-based processing. For example, they propose taking 27 packets with 9 features each and creating a $9 \times 9 \times 3$ image, where three channels RGB (red, green, blue) are provided for a single pixel. This is then used as input into their CNN architecture.

Gao [108] proposed a CNN-LSTM architecture which reduces false positive rates in their network intrusion detection system for cloud computing environments. Numerical and character attributes from the KDD CUP 99 test data set are first pre-processed into a 10×42 matrix, where window size $w = 10$ and the number of features in the data set is 42. The data set is transposed to an image dataset where each matrix W is mapped to a 10×42 grayscale image. If v_{ij} represents the value of a pixel in matrix/image W at position i, j , the values $v_{ij} \in [0, 255]$. To reduce the data interval and normalize values, the following formula is used:

$$x = \frac{x - \bar{x}}{\delta} \quad (5.5)$$

where x is the current value of a certain data in the sequence, \bar{x} is the mean value of the data in the series, and δ is the standard deviation of the sequence. Finally, this system adds an additional attention mechanism layer to improve classification accuracy:

$$\begin{aligned} u_t &= \tanh(W_w P_t + b_w), \\ a_t &= \text{softmax}(u_t^T, u_w), \\ v &= \sum a_t P - T, \end{aligned} \quad (5.6)$$

where u_t is the attribute representation of P_t , W_w is the context vector randomly generated during training, a_t is the importance weight, and v is the high-level repre-

sensation obtained through importance weight summation on P_t .

Jo et al [3] use a CNN-LSTM hybrid model to perform network intrusion detection, but first pre-process packets into image representations. Namely, they adapt the NSL-KDD dataset which has 41 fields into pixel representations. Continuous fields are first normalized to a $[0 - 255]$ range to create a grayscale image pixel value. Symbolic fields such as “protocol type” are represented by the same kind of pixel value, but normalized based on the number of possibilities. For example, a symbolic value with three possibilities will be normalized as 0, 128, or 255 [3]. 28×28 images are created and processed with a 5×5 kernel as shown in Figure 5.1.

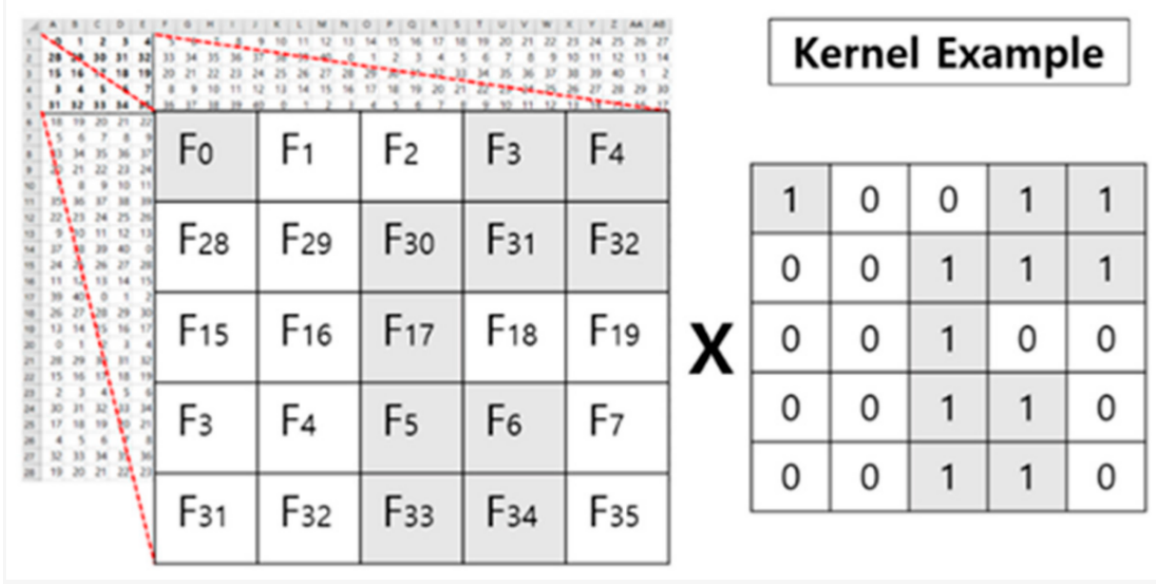


Figure 5.1: Example application of kernel to pixelated packet data [3]

Deep Packet [38] is a well-cited work in deep learning for traffic classification. In this system, a stacked auto-encoder and CNN are employed to classify traffic by type and application as labeled in the VPN/non-VPN UNB ISCX dataset [90]. Their solution achieves 94% accuracy in the traffic classification task per traffic type, and 97% accuracy in the application classification task. Incorporating the same dataset as well as additional, more recent flows, we are able to achieve higher accuracy in

RTP detection with a less computationally complex framework in MAPLE.

5.3 Methodology

We expect some level of standardization (i.e. method names or status codes) or some commonalities such as user names, device details, semantic contexts, and other conversation or implementation-specific indicators even in highly entropic or variable packet data. Using this hypothesis, MAPLE transforms packets into grayscale images which will appear similar to one another. Thus, machine learning algorithms designed for image analysis may record the hidden features generated from the similarities between the two images. We base this assumption on previous work [106] which has shown similarities in grayscale images rendered from packet data in order to classify traffic by application and protocol type which are detectable by CNNs. Figure 5.2 asserts the validity of this assumption by showing comparative images rendered from packets in our combined dataset previously used and explained in detail in chapter 4 of this work.

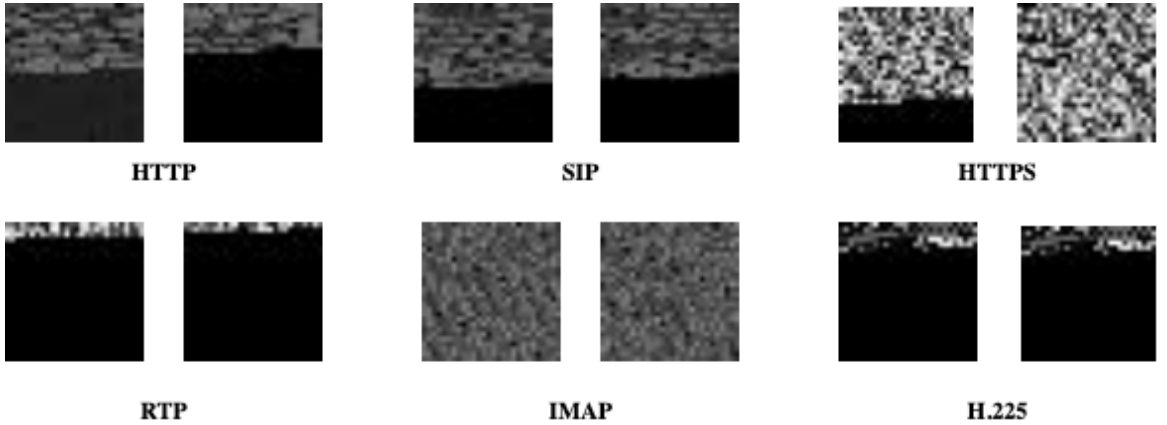


Figure 5.2: Grayscale images generated from packets in our combined dataset used in the following experiments

To shape packet payloads into an image, we extract the payloads and convert them from byte encoding to a normalized decimal integer $i \in [0, 255]$ [3]. For image size, we chose a 28×28 representation, for a total of 784 input features. Padding is applied

in the cases where the payload is shorter than 784 bytes, and is otherwise truncated.

5.3.1 LeNet Model Configuration

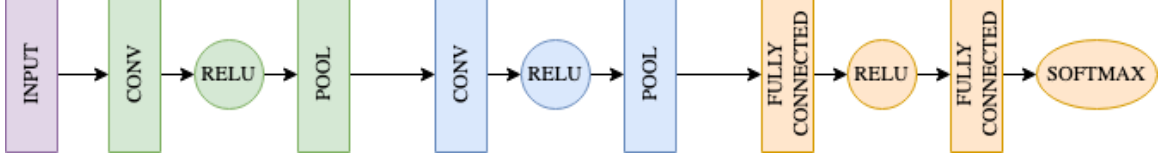


Figure 5.3: The architecture of our LeNet models

LeNet was one of the standard models we selected in order to test convolutional neural networks against the RTP detection problem. LeNet has a low complexity, so has higher potential for practical use in real-time, line-rate packet inspection systems than deeper models which require more time. Figure 5.3 shows the set up of the model which we use as a baseline. We designed two separate models where one employed 6 and 16 kernels per convolutional layer (A) and another that used a 16 and 32 kernels (B). Each model contained two fully-connected, dense layers of size 1024 with a binary softmax classifier at the output layer.

5.3.2 ResNet Model Configuration

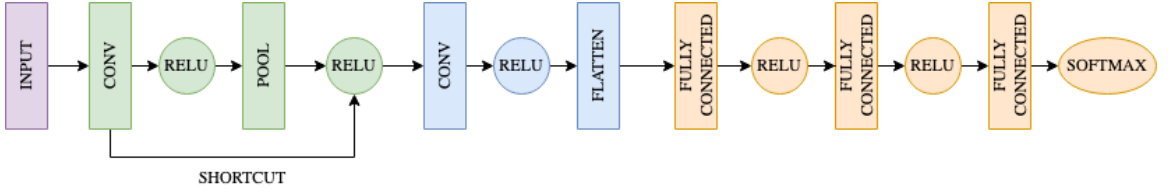


Figure 5.4: The architecture of our ResNet models

If they are to be deployed at scale, artificial intelligence solutions must be capable of line-rate processing while still being able to perform the classification task to the required level of accuracy. This can be difficult to quantify as the definition of line-rate varies per network environment. Still, we assert that the identity mappings of ResNet do not introduce additional complexity [105]. Thus, residual mapping is introduced as a potential solution to adding additional layers of representation and thus improve

classification through more hidden features, while also minimizing overhead.

We developed three ResNet-based models for the MAPLE system’s experiments. The first model (C) consists of three convolutional layers with 32, 32, and 64 kernels of dimension 3×3 , and three dense layers with output sizes of 64, 32, and 2 respectively. The second model (D) follows a similar architecture except the values are halved. Convolutional layers had 16, 16, and 32 kernels per layer and the three dense layers had output sizes of 32, 16, and 2. Model (E) corresponds to the same number of kernels and output sizes as (C), except it uses a kernel dimension of 7×7 . The ReLU activation function and adam optimizer was used in all the model configurations, and a final layer used softmax for normalization and classification. The loss function used for training was categorical cross-entropy, and we also employed a dropout layer to reduce overfitting.

Table 5.1: Configurations of each model used in the experiments

Type	Model	# Kernels	Kernel Size	FC Dim
LeNet	A	3×3	6, 16	1024, 2
	B	3×3	16, 32	1024, 2
ResNet	C	3×3	32, 32, 64	64, 32, 2
	D	3×3	16, 16, 32	32, 16, 2
	E	7×7	32, 32, 64	64, 32, 2

5.4 Experiments and Results

We ran several experiments to test the model’s ability for binary RTP/non-RTP detection as well as multi-class protocol identification. All tests were performed on a single CPU of a 1.6 GHz dual-core Intel i5 processor with 16 GB DDR3 RAM. In the binary confusion matrix, true positive (TP) indicates correct classification of data. True negative (TN) is correct classification of data as non-RTP. false negative (FN) implies incorrect identification of traffic as non-RTP, and false positive (FP) is the incorrect classification of data as RTP. We use measurements of precision, recall, and F1-score for model evaluation, defined as follows:

$$\begin{aligned}
Recall &= \frac{TP}{TP + FP} & Precision &= \frac{TP}{TP + FN} \\
F1\ Score &= \frac{2(P \times R)}{P + R}
\end{aligned} \tag{5.7}$$

5.4.1 RTP Detection

For the binary RTP/non-RTP classification, we merged SRTP and RTP traffic into an RTP label and re-labeled all non-RTP and non-SRTP traffic as non-RTP. For testing the different configurations of the MAPLE model, we first processed traffic and transformed it into a matrix image, and then used it as input for each of the models. We performed random under-sampling to avoid bias in the dataset, and then split the data into 60%/40% for training and testing. We ran the multi-model using ALPINE and PALM with the same datasets with the results in Table 6.1 to establish a baseline for comparison. Then we ran tests with training for 3 epochs for each CNN model. As we plan to deploy this technology in a real RTP detection and deep packet inspection system, we also measured classification throughput in order to determine how much traffic the MAPLE system would be able to process given the detection model configuration. Results are provided in Table 5.3.

Table 5.2: Classification results of RTP vs non-RTP detection for ALPINE and PALM

Class	P	R	F1
<i>Non-RTP</i>	0.73	0.77	0.75
<i>RTP</i>	0.76	0.71	0.73

Generally, the ResNet model performed exceedingly well at the RTP detection task, exceeding ninety-nine percent accuracy across all the configurations. The LeNet models performed marginally worse in terms of accuracy, but the smaller of the two models had the highest throughput of any of the tested configurations. Another configuration choice which may have significant impact on model accuracy is the

Table 5.3: Classification results of RTP vs non-RTP detection for all MAPLE models

Model	P	R	F1
A			
<i>Non-RTP</i>	0.96	0.99	0.98
<i>RTP</i>	0.99	0.96	0.98
B			
<i>Non-RTP</i>	0.97	0.99	0.98
<i>RTP</i>	0.99	0.97	0.98
C			
<i>Non-RTP</i>	1.00	0.99	1.00
<i>RTP</i>	0.99	1.00	1.00
D			
<i>Non-RTP</i>	1.00	1.00	1.00
<i>RTP</i>	1.00	1.00	1.00
E			
<i>Non-RTP</i>	1.00	1.00	1.00
<i>RTP</i>	1.00	1.00	1.00

Table 5.4: Throughput results of RTP vs non-RTP detection for all MAPLE models

Model	Accuracy	Mb/s
A	0.975348	12.5
B	0.978258	9.28
C	0.996934	7.07
D	0.995908	10.41
E	0.996847	5.3

number of epochs trained. Real systems place an emphasis on minimizing down time of a system, but in actual deployment many models may be trained offline and then embedded. In systems where training time is not a factor of performance, we can thus increase the number of epochs with little real impact. In order to test the efficacy of such a design choice, we re-ran MAPLE models A and C with 10 epochs to see if performance improved. Results in Tables 5.5 and 5.6 show that ResNet appears to learn faster than LeNet, as it gained very little improvement with additional training time. Interestingly, with additional training time LeNet approaches similar accuracy to ResNet but still maintains higher throughput.

Table 5.5: Classification results of RTP vs non-RTP detection for MAPLE models A and C with additional training time

Model	P	R	F1
A			
<i>Non-RTP</i>	0.99	0.99	0.99
<i>RTP</i>	0.99	0.99	0.99
C			
<i>Non-RTP</i>	1.00	1.00	1.00
<i>RTP</i>	1.00	1.00	1.00

Table 5.6: Throughput results of RTP vs non-RTP detection for MAPLE models A and C with additional training time

Model	Accuracy	Mb/s
A	0.990499	14.99
C	0.998172	7.38

5.4.2 Protocol Auto Detection

As previously mentioned, many middlebox technologies are interested in multi-classification problems where many different applications or protocols need to be identified. Thus, we expand our experiments to a 26-protocol problem using the combined dataset described in chapter 4. We again use models A and C, LeNet and ResNet respectively, for this expansion. As this is a much harder problem, we expect to need more training time. We ran each model with 10 epochs and 50 epochs, and recorded the macro averages in Table 5.7.

Table 5.7: Results of multi-class detection for models A and C as macro averages

Model	P	R	F1	Accuracy
A - 10 epochs	0.78	0.78	0.77	0.780920
50 epochs	0.83	0.83	0.83	0.832446
C - 10 epochs	0.84	0.84	0.83	0.835993
50 epochs	0.86	0.85	0.85	0.849528

Increasing the training time for the LeNet model showed significant improvement in accuracy, but not likely to a level needed for deployed systems so there is a need to build upon the baseline. The ResNet model approaches better accuracy with more

training time.

5.4.3 User Fingerprinting

In order to test capability to track user streams, we isolated a source/destination IP pair in the traffic and created a separate user label for it. In this way, we associate a particular traffic stream with a person of interest. Then, we configured MAPLE and a multimodal version of ALPINE to use no network layer features, and instead rely on packet length, flags, and analysis of payload contents and spatial features to make the classification. Results in Figure 5.5 show this approach as a promising solution for tracking user streams across changing network configurations.

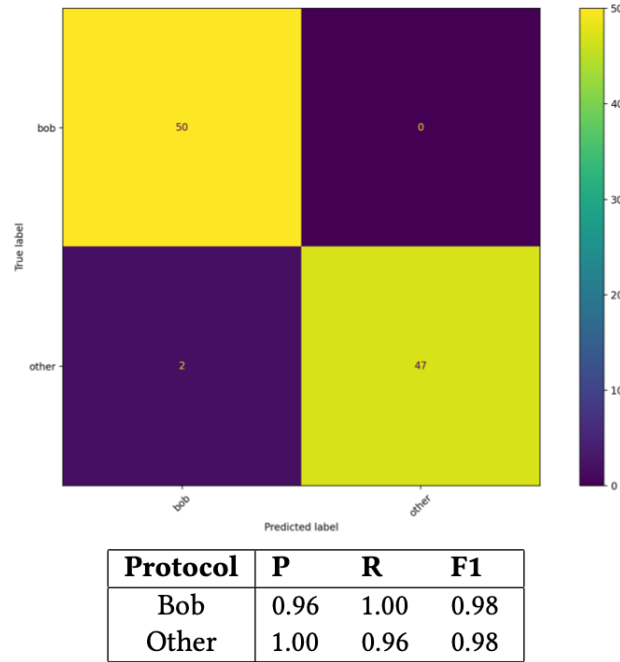


Figure 5.5: A confusion matrix of tracking user Bob across traffic streams in a changing network environment

5.4.4 Model Performance and Throughput

Decreasing the number of kernels may have slightly impacted accuracy, but improved overall throughput significantly and lessens memory footprint. In a real network environment, it may be worth consideration to sacrifice some accuracy in order

to be able to monitor more traffic or more effectively load-balance the received input depending on hardware capability. Thus, the ideal model configuration is often environment-dependent and inspired us to incorporate configuration capability in MAPLE so that the deployed system may best suit the environment. We provide the throughput of each model configuration in our experiments in Table 5.4.

5.5 Discussion and Limitations

MAPLE works well on input data such as a payload which may be distinct from one data sample from another, but can be divided and matricized into units for comparison (i.e. turned into a grayscale image of uniform dimension). For deep packet inspection, this model is able to create latent representations of payload data which uniquely identify traffic of different types at a higher dimension than is considered by current signature-matching or filter-based solutions used in industry. Our implementation for this initial deployment solves the RTP detection problem with high accuracy using a minimal framework ideal for line-rate. The system could be further optimized by running these models in parallel; as the solution relies on per-packet analysis, throughput could be increased across models. One limitation of the proposed encoding is that individual packets must be large enough to create a significant image. Thus, this model would not work well on short packets or packets with no payload data. In this case, we propose using a combination of header features and payload when available would be a wiser strategy.

CHAPTER 6: DATE - POINT CLOUD REPRESENTATION AND DENSITY CLUSTERING ANALYSIS OF PACKETS

6.1 Introduction

While image generation and matrix representation expand packet data to two-dimensional space, we wanted to explore other geometric representations of data and deepen the dimensionality. In geometric topology, point clouds are another way to generate data shapes for anomaly detection and comparison. For this purpose, we designed DATE, a cluster **D**ensity **A**nalysis-based **T**ensor **E**ncoder model which expands network packets into three dimensional point clouds, which we then extract features of regarding cluster density in order to find commonalities and classify based on data shape. This spatial expansion is a novel application of latent representation learning to network traffic and the problem of RTP detection or detecting weakly signaturred traffic and could expand to other protocol problems.

6.2 Previous Work

There is little previously published work in generating three-dimensional point cloud representations of packets, making DATE a novel contribution to the field of packet processing. LiDAR (light detection and radiation) systems generate packet data which has been transformed into both two-dimensional matrices and point clouds [109]. This work introduces the problem of spatial correlation in packet data, namely that packets in their raw state are not usually uniform or in a state which may be processed using the spatial or geometric strategies employed by image and computer vision algorithms. Therefore, engineers must perform data compression or pre-processing to create the necessary uniformity.

Costeux et al [110] worked on the fast detection of Skype and other RTP-based telephony traffic. They specify several fields from the RTP and encapsulation header which can be used to filter out non-RTP traffic. We employ this technique in the middlebox as an initial filtering step in our end-to-end data processing pipeline. Kmet et al [99] build on Costeux by incorporating additional header features for per-packet selection, and adding a flow-based solution to reduce the over-selection problem. They were able to minimize mis-classification to nearly zero by buffering up to 10 packets of the RTP stream. The synchronization source number from the header is used along with timestamps to check for proper increment and stream correlation. We focus on per-packet identification only.

Some early research [111] describes expanding NetFlow data generated by routers into manifolds. We reference this as an initial introduction of spatial expansion of traffic data in order to perform dimensionality reduction, which is a key theoretical linkage to our work.

6.3 Methodology

DATE describes the density-based clustering analysis which is a novel contribution of this work. This process attempts to map data to a three-dimensional space, then use spatial features to generate hidden features.

6.3.1 Point Cloud Creation

Point clouds are a set of Cartesian coordinates (x, y, z) which represent some points in a three-dimensional space. Point clouds may be used in computer vision software as a method for mapping high-dimensional shapes into lower-dimensional space [112]. As a latent representation, this method can be used temporally to express the movement of objects in three-dimensional space. An airplane can be modeled as a three-dimensional point cloud, and its movement described as a series of clouds c_0, c_1, \dots, c_t where t represents a number of time intervals and c_n the point cloud generated at

a given time slice. Like the work of Quach et al [113], we do not combine multiple representations in a time series but instead analyze the static representations much analagous to manifolds in three-dimensional space. We use the geometrical representation captured here for cluster analysis.

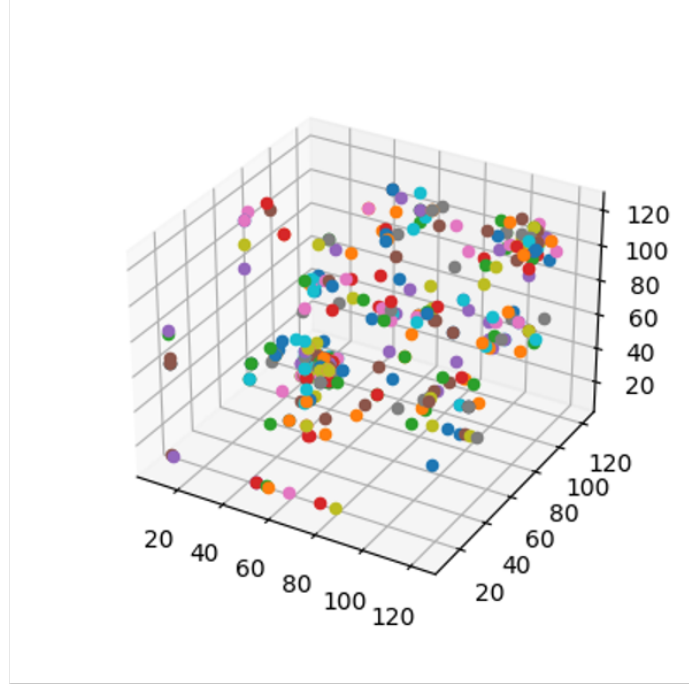


Figure 6.1: 3D point cloud made out of a Session Initiation Protocol (SIP) packet

In the DATE model, payload data is extracted from the packets and truncated to 1024 bytes. If the payload is shorter than this, we add padding for normalization. Byte values v are converted to decimal values $v' \in [0, 255]$. For example, the byte string, '0x68', '0x65', '0x6c', '0x6c', '0x6f', would first be transformed to decimal values (104, 101, 108, 108, 111). Then, the sliding window would produce the points corresponding to coordinates:

$$\begin{aligned}
 C = \{ & (104, 101, 108), \\
 & (101, 108, 108), \\
 & (108, 108, 111) \}.
 \end{aligned} \tag{6.1}$$

The coordinates are mapped to a three-dimensional space to form point clouds like Figure 6.1. We expect that similar point clouds will be generated from packets with similar data.

6.3.2 Density-Based Spatial Clustering of Applications With Noise (DBSCAN)

Density-Based Spatial Clustering of Applications with Noise (DBSCAN) uses three different point classifications. *Core points* make up the centers of clusters. Individual core points are labeled based on the number of neighbors they have compared to their neighbors, or their degree of importance. *Border points* are those which have the least number of neighbors but are still comparatively relevant due to being within the range of a core point. These comprise the edges of the point cloud. *Outlier points* are neither borders nor cores, do not belong to a cluster, and are considered noise [114]. The DBSCAN algorithm visits each point and classifies them as such in order to paint a picture of where clusters of data are as in Figure 6.2.

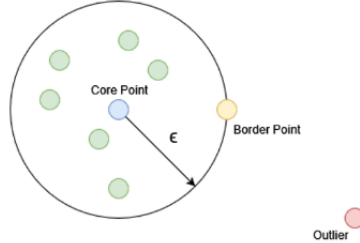


Figure 6.2: Illustration of types of points in DBSCAN

We first create the point clouds and then process them using DBSCAN in order to derive the natural clusters that exist in the data. DBSCAN groups points that are close together utilizing a Euclidean distance for measurements. The value of the Euclidean distance that determines proximity and the minimum number of points needed to form a cluster are configurable parameters: ϵ and a minimum number of samples. The ϵ value is the radius for each circle drawn around each point to query density, and the minimum number of samples is the minimum number of points required to be inside the drawn circle for it to be considered a Core point.

Generally, the number of samples is twice the dimensions, but it must be at least the number of dimensions. If it is set below this threshold, then singular points or two points that are close together would constitute a cluster, which would not be an accurate way of determining where data is clumped [115]. In our experiments, we configured DBSCAN to perform with $\epsilon = 7$ and `min_samples = 7`. This would be enough to create more clusters in a less populated point cloud. An ϵ value of 7 also allowed for a greater number of recognizable clusters [116].

Like many machine learning models, the parameters used to tune DBSCAN greatly affect performance rates. It is possible to create more clearly defined point clouds within the three-dimensional space through this parameter adjustments. The ϵ value needs to be set based on general rules. If the value of ϵ chosen is too small then too many clusters will be created which over-emphasize the importance of noisy points. However, if set too large, then clusters will merge together and shapes become less distinct. One way that a value of ϵ can be estimated is by examining the input data and finding the average distance between each point and its K -nearest neighbors. Choosing an exact value for ϵ is a difficult problem for packet data which is highly variable. The ϵ value must be estimated from a conglomeration of the packet's K -nearest neighbors where K is equal to `min_samples`.

Because each packet is highly variable, choosing an exact value for ϵ is a difficult problem. An estimation of ϵ must be taken from a conglomeration of the packet's k -nearest neighbors data where k is the value of `min_samples` chosen. As shown in the figure above, the point of each line with the greatest curvature is considered the ideal value for ϵ . A line has been illustrated in Figure 6.3 where an estimated average is made for every sample. The variation in an ideal ϵ value could be a major factor in why some packets are recognizable in their cloud state and others are not [116].

We extract the following features as the feature vector for classification from the DBSCAN results:

- clusterCount = number of clusters
- averageClusterSize = average cluster size
- standardDeviation = standard deviation
- noisePercent = percent of cloud containing noise

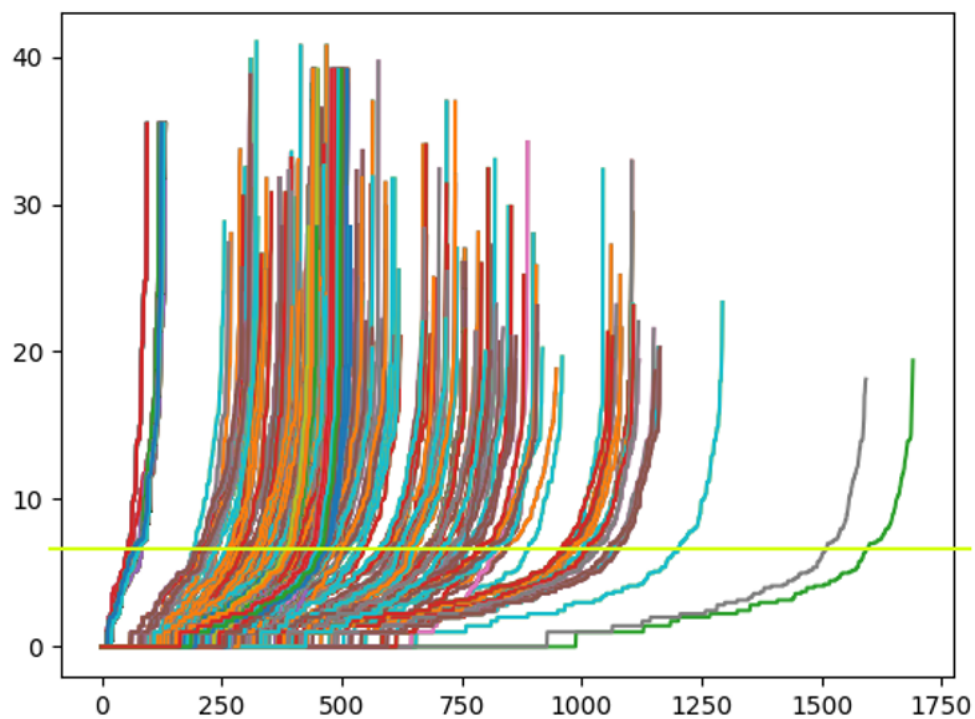


Figure 6.3: K-nearest neighbors graph generated for packets in the dataset using a `min_samples = 6`

6.3.3 Packet Classification

For deep learning, we feed the extracted cluster features forward into a two-layer multilayer perceptron unit (MLP) with a final softmax layer for classification. We use binary cross entropy as the loss function.

6.4 Experiments and Results

We experimented with the RTP detection problem that was previously attempted with the MAPLE model in chapter 5 in order to assess the ability to detect weakly signed data. Then, we extend this to H.225 detection, which is another data transfer protocol with a weak signature. Lastly, the model performance and throughput is assessed for its application to deep packet inspection and deep learning at scale.

All tests were performed on a single CPU of a 1.6 GHz dual-core Intel i5 processor with 16 GB DDR3 RAM. In the binary confusion matrix, true positive (TP) indicates correct classification of data as RTP. True negative (TN) is correct classification of data as non-RTP. false negative (FN) implies incorrect identification of traffic as non-RTP, and false positive (FP) is the incorrect classification of data as RTP. We use measurements of precision, recall, and F1-score for model evaluation, defined as follows:

$$\begin{aligned} Recall &= \frac{TP}{TP + FP} & Precision &= \frac{TP}{TP + FN} \\ F1\ Score &= \frac{2(P \times R)}{P + R} \end{aligned} \tag{6.2}$$

All tests were performed with undersampling to appropriately balance the dataset, and a 60%/40% training and test split. Epochs represent the number of passes made over the training data in order for the neural network to fine tune itself under supervision. Although more epochs do not guarantee better performance, it is often the case that a larger number of epochs within reason can improve metric results. However, too many epochs can cause overfitting to the training data and hurt generalizability.

6.4.1 RTP Detection

Our primary goal in designing DATE was to capture spatial hidden representations which may be present in these geometric reinterpretations of packets. These features may present advantage to our classifier as dynamic ports and weak signatures found in

these protocols make traditional signature-based and port-based methods ineffective. The classification results of the DATE model are provided in Table 6.2. Table 6.1 show the combined ALPINE and PALM technique’s results for the RTP classification task. There is clear improvement for the DATE model.

Table 6.1: Classification results of RTP vs non-RTP detection for ALPINE and PALM

Class	P	R	F1
<i>Non-RTP</i>	0.73	0.77	0.75
<i>RTP</i>	0.76	0.71	0.73

Table 6.2: Classification results of RTP vs non-RTP detection for the DATE model

Class	P	R	F1
1 Epoch			
<i>Non-RTP</i>	0.85	0.84	0.84
<i>RTP</i>	0.83	0.85	0.84
10 Epochs			
<i>Non-RTP</i>	0.94	0.81	0.87
<i>RTP</i>	0.83	0.95	0.89
20 Epochs			
<i>Non-RTP</i>	0.96	0.79	0.87
<i>RTP</i>	0.82	0.96	0.89

6.4.2 H.225 Detection

H.225 handles the registration and call setup for certain VoIP architectures using the H.323 protocol suite. Similar to RTP, H.225 data can arrive at the endpoint first before H.323 which it must be paired with. It also does not have a strong signature. Thus, we additionally tested DATE’s ability to detect H.225 traffic from non-H.225 traffic. Our setup of the dataset labels and the DBSCAN configuration was similar to the RTP test and results are given in Table 6.4. The performance of ALPINE and PALM are provided in Table 6.3 for a baseline.

6.4.3 Model Performance and Throughput

In order to benchmark the system for real deployment, we recorded the time classification took. While this is heavily system dependent and showed some variance

Table 6.3: Classification results of H.225 vs non-H.225 detection for ALPINE and PALM

Class	P	R	F1
<i>Non-H.225</i>	0.62	0.64	0.63
<i>H.225</i>	0.63	0.61	0.62

Table 6.4: Classification results of H.225 vs non-H.225 detection for the DATE model

Class	P	R	F1
1 Epoch			
<i>Non-H.225</i>	0.98	0.83	0.90
<i>H.225</i>	0.85	0.99	0.92
10 Epochs			
<i>Non-H.225</i>	0.91	0.87	0.89
<i>H.225</i>	0.88	0.92	0.90
20 Epochs			
<i>Non-H.225</i>	0.98	0.83	0.90
<i>H.225</i>	0.85	0.99	0.92

across parameter tuning or model configuration, the average classification time for a single packet by DATE was 0.15823 seconds.

6.5 Discussion and Limitations

Of the models tested so far in this work, DATE is the most computationally complex with the lowest throughput. In future work, we propose implementing multi-threading and multi-core processing as potential optimizations. We observed that point cloud generation was a pain point for the system in terms of cycles, and consider offloading such repetitive calculations to a specialized hardware such as FPGA [117] when available in the deployed system. Because DATE is able to normalize data and map it into a three-dimensional cluster space, it would be well-suited to heterogeneous data such as sensor data combined with packet data from an Internet of Things (IoT) device. This would create a multi-embedding capable of correlating these data inputs for machine learning tasks. Simpler problems may be solved with less complex approaches such as regular expression matching on plaintext SIP traffic, or using the locality-sensitive hashing strategy toward device fingerprinting.

We propose both MAPLE and DATE as potential methods for generating hidden, latent-space representations of traffic as their capabilities extend to different problem areas. MAPLE works well on input data such as a payload which may be distinct one data sample from another, but can be divided and matricized into units for comparison (i.e. turned into a grayscale image of uniform dimension). On the other hand, DATE has the potential to expand to include other non-network features and represent more heterogenous data. For example, input sensor data from IoT devices or light detection and ranging (LiDAR) equipment have used DBSCAN for data processing and normalization [118]; it could be combined with cloud/network data inputs as an additional embedding model for classification problems. While for the RTP problem MAPLE provides high accuracy with more efficiency than DATE, there is a trade-off within the embedding space as DATE may be able to represent data of higher complexity in other problems for future work.

CHAPTER 7: FORAGER - TOWARDS MULTI-EMBEDDINGS AND ENSEMBLE CLASSIFICATION

7.1 Introduction

In the previous work with ALPINE and PALM we hypothesized that enabling both models to run together and combine their votes would yield the most accurate and agreed-upon classification results. The results of the multi-embedding experiments in this work proved this to be more accurate. Multiple classifier systems (MCS) are ensembles and are widely used in pattern recognition applications and recognized as more effective than any single embedding or classifier. There are at least three reasons identified by Dietterich as to why multi-classifier approaches perform well [119]:

- Statistically, there is a set of H hypotheses for which the learning algorithm is performing a search for the best hypothesis. The risk of choosing the wrong classifier is reduced when using multiple for unseen data.
- Computationally, optimal training is an NP-hard problem. Because the optimum depends on the starting point, the optimal learning algorithm may be different. By using an ensemble, a better approximation of the unknown function may be obtained.
- Representationally, the weighted sums of hypotheses from H allows for the expanding of space of representable functions. Simply put, the more embeddings we have, the more features, and the more ways to represent and compare data.

In addition to multiple classifiers, in the workflow process using multiple types of embeddings has proven more effective in several recent, real systems than any single approach [120, 121, 122, 107, 123, 124, 125, 126, 127].

As we have seen in previous chapters, each model unveils different hidden representations which are appropriate for different classification problems and contexts. In the real network environment, these classification problems are not isolated. Traffic may need to be distinguished between encrypted and plaintext, for example, and then protocols identified. Then, encrypted traffic may be profiled, or the plaintext traffic characterized by embedding for user fingerprints. Thus, a combined approach or toolkit is apt for the many kinds of problems real world network analysts face.

To meet this need, we propose FORAGER, a combination of data mining and hidden representation learning approaches to deep packet inspection and network traffic classification. This highly configurable toolkit includes the ability to extract and transform header features into locality-sensitive hashes indexed in an LSHForest using the ALPINE technique. Payloads may be combinatorially analyzed using either the PALM, MAPLE, or DATE approaches, or a combination of votes from these models. By considering the data from multiple angles, we achieve results of higher accuracy and a much greater understanding of the network traffic than other state-of-the-art systems.

To illustrate the capability of FORAGER and its diversity of traffic applications, we performed a case study on profiling the traffic on port 443. Due to the rise of encryption and its standard usage for Hypertext Transfer Protocol Secure (HTTPS) encrypted traffic, port 443 is often left un-monitored by packet inspection on large-scale surveillance systems or default-allowed by firewall software. While the majority of Internet traffic is now safely secured behind transport layer security, there are threat actors who have found ways to utilize port 443 as a covert data channel for surveillance and firewall evasion, malware transport, data exfiltration, and nefarious activity under the guise of Internet anonymity. For increased security and surveillance, it is necessary to be able to filter legitimate, HTTPS traffic from illicit activity on this port. We propose using FORAGER, our network traffic classification and

profiling toolkit, to dissect this traffic stream for hidden content. Using FORAGER and its individual models, we are able to extract plaintext and compressed traffic from encrypted streams with high accuracy, profile Tor-based, VPN-based, and encrypted traffic streams by application and data type, and identify malicious instances of particular protocols designated to run over HTTPS. Our models combine several approaches to hidden feature extraction from packets via spatial representation learning to precisely and efficiently classify packets, fortifying network security while still preserving user privacy through this crucial gateway.

7.2 Previous Work

In our threat model and in the works we compare against [6, 7, 38], we choose to prioritize making classification using only a single packet from any point in the overall traffic flow. This contrasts systems which require entire flows in order to classify [37, 38, 39, 40, 41].

7.2.1 VPN Traffic Profiling

Deep Packet [38] uses a CNN to perform traffic classification into traffic types and application types as labeled in the VPN/non-VPN UNB ISCX dataset [90]. Their solution achieves 94% accuracy classifying traffic by type, and 97% accuracy in classifying by application. Like our work, they use a single packet in the classification task, and do not rely on flow-based features. We note that their work does not perform as well in Tor traffic classification, leading us to conclude that Deep Packet may not generalize as well as FORAGER. Cui et al use only header information from the ISCX datasets to classify traffic over flows with 99% accuracy [35]. While a novel solution with enhanced privacy as it considers only header features, it still requires flow-based features of the header and therefore multiple packets. Zou et al use the spatial features of the first packet of a flow extracted and applied to a CNN and analyze the time series features of the next three packets using a long short-term memory

network (LSTM) [36]. This method shows high accuracy on the dataset but requires flow-based information.

7.2.2 Darknet Detection and Profiling

Lotfollahi et al also study the ISCX dataset for Tor data and profile again by application and traffic type using a single packet approach [38]. Tor traffic profiling has been further explored by Lashkari et al [7] in their DIDarknet system using a single packet, where they also proved accurate results using the ISCX Tor dataset. Other works have used flow-based features for both detecting Darknet activity (binary Tor/non-Tor classification) and profiling Tor traffic into specific applications or traffic types (e.g. streaming, chat, web browsing). Iliadis et al [128] use an amalgamated version of the ISCX Tor and VPN datasets, re-labeled CICDarknet2020, in order to create a binary classification of Darknet and regular traffic, and a multi-class problem of Tor, non-Tor, VPN, and non-VPN. We perform a similar analysis in our experiments and results of the same data. The work uses all 85 features published in the dataset, including flow-based features, and the full traffic information. Using a set of machine learning algorithms, they are able to achieve up to 98% accuracy; however, in comparison our system uses only a single packet and a minimal set of features to achieve similar results. Ma et al [129] also achieved noted success in the community classifying the same dataset using a CNN with Root Mean Square (RMS) propagation. They also require flow-based features of the dataset, and would therefore require multiple packets and flow information in a real system. They also only consider Darknet detection without additional profiling and achieve 95% accuracy, lower than ours in the same task. Sarkar et al [130] introduce a Generative Adversarial Network (GAN) in order to stabilize their deep learning model across the same ISCX Tor dataset, and achieve approximately 98% accuracy, as well. Using flow-based features, they are able to provide a 95% accuracy in the profiling task across the eight described traffic classes in the dataset; however, they do not provide the results broken down

individually so it is difficult to conclude if the model performs evenly well across classes or not. Choorod et al [6] use a statistical analysis solely of the extracted features of the encrypted packet payload. They find that certain sizes are common across traffic types, as well as certain regions of the encrypted payload. They use Charcount from the Posit Text Profiling Toolset [131] as a ratio value for input to their decision tree-based algorithm. They achieve highly accurate results we compare against in the Experiments section of this work.

7.2.3 Filtering Compressed/Plaintext Traffic

Plaintext filtering can be done through statistical entropy tests [19] and has been used to find private medical or geo-location data [19, 132], covert malware [133], and general unencrypted traffic like file transfers or emails. Modern compression and encryption algorithms present such similar high entropy that distinguishing between them is a difficult problem. The HEDGE [134] system is able to parse encrypted from compressed traffic with between 60-94% accuracy using the Chi Square and a subset of the NIST SP 800-22 tests dependent upon the size of the packet used. They compare against Hahn et al, who achieved a maximum of 66% accuracy using machine learning models [135].

7.2.4 Intrusion Detection over HTTPS

Zain ul Abideen et al [136] consider the problem of detecting and profiling VPN traffic on port 443, where their system uses non-encrypted, packet level features to determine several VPN types in the traffic flow. They distinguish between VPN configurations, including Tor. Furthermore, they employ a technique using DNS queries to identify malicious or illegitimate VPN servers. This is an interesting notion which may be challenged by the increased use of encrypted DNS and techniques like DNS-over-HTTPS. Identification of SSH flows in encrypted traffic has been accomplished with success when using flow data [?]. Detecting brute-force and dictionary attacks

over SSH is also a research problem area from which our work on port 443 is inspired [137].

7.3 Methodology

Details of the network traffic classification systems used in FORAGER can be found in previous chapters of this work. In addition to technical methods, we provide practical application through analysis of threat models for which FORAGER is useful.

7.3.1 Threat Model

To bring the problem to life, we envision a criminal case in our *threat model*. Bob is a member of an organization embezzling funds through an international shipping company. There are several nefarious tasks which Bob needs to accomplish that he could utilize port 443 exploitations to achieve:

- Create a *covert data channel* in order to pass along sensitive information like receipts and records while avoiding detection,
- *Communicate anonymously* with criminal network colleagues to avoid tracing of colleagues, his own geolocation data, and evidence,
- Establish *remote access* to a machine from endpoints configured to disallow ports like 22,
- *Exfiltrate data* from a secure source.

Alice is an intelligence officer who is monitoring Bob's network data with reasonable suspicion that he is engaged in criminal acts. We will discuss vulnerabilities in Alice's current surveillance framework and how to use FORAGER to capture and record evidence which can be used to thwart Bob's plans and provide legal support in the case against his embezzlement and fraudulent company.

7.3.1.1 Real-Time Monitoring

In our scenario, Alice possesses a middlebox technology configured to monitor traffic on port 443. It is common for routers or middleboxes to default-ignore traffic on port 443 as there is usually an overwhelming amount of encrypted traffic on this port that is unprocessable [138]. Attackers who are aware of this strategy may send unencrypted or sensitive information over port 443 in order to avoid surveillance and detection. In criminal operations, these packets could contain significant data. Thus, it is important for middlebox technologies to have a methodology for scanning traffic on port 443 to determine if the data is truly encrypted. If the data is plaintext, it may be directly processed. If it is highly entropic but only compressed, follow-on processing may be capable of decompressing the data. Currently, Alice configures her middlebox to default-ignore traffic on port 443, instead only monitoring port 143 for unencrypted emails. In order to avoid detection by possible surveillance technology, Bob routes his emails through port 443. The receiver is aware of this and listens for this traffic with the IMAP server on the configured port so that it is received. Thus, Bob is able to transmit content in a way that is receivable by the partner as email/IMAP at no additional security overhead but is undetected as Alice's middlebox does not scan port 443.

7.3.1.2 Tunnel Tracking

Tor allows anonymous Internet space for pervasive action and may be used to access content which is blocked by firewalls, whether it be censored or illegal activity. While possibly encrypted and encapsulated with layers of addressing in order to anonymize routing, traffic may be synchronized by timing and further analysis in order to determine endpoints [139]. Thus, tracking Tor and being able to further profile Tor usage and content is a useful operation for intelligence and law enforcement. Network Address Translation (NAT) in VPNs prevents traditional ISP tracking. The ability to

detect VPN traffic streams and profile their traffic and applications can be additionally useful for network monitoring. In the scenario, Bob is using a VPN to connect to a forum website. He wishes to remain anonymous during this communication. Alice is following Bob's messages, and sees that he regularly communicates using HTTP and IRC protocols on this platform. She configures her middlebox device to promote traffic on several five-tuple-based traffic flows which she knows to be currently associated with Bob's devices. However, dynamic re-configuration of the tunneled network layers causes Alice to lose the data stream and she is no longer able to monitor his messages and record the activity.

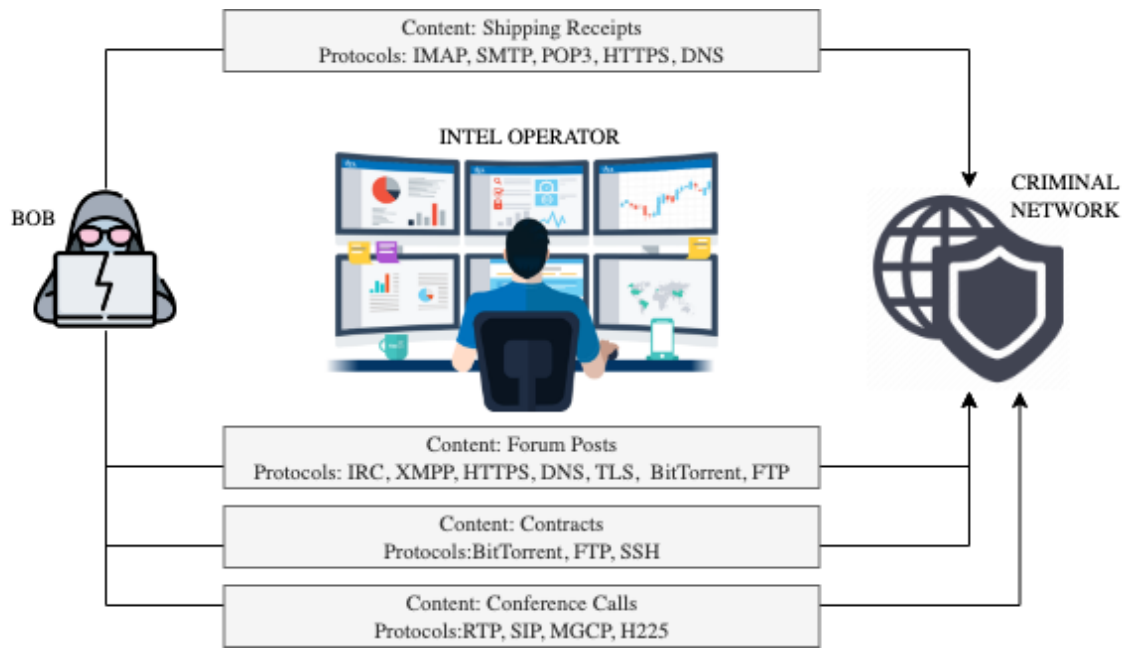


Figure 7.1: A depiction of intelligence operations intercepting various types of traffic Bob is transmitting in connection to his extended criminal network

7.3.1.3 Covert Data Channels In

Domain name security (DNS) has long been exploited by attackers due to many security loopholes and its historical transmission in plaintext. Malicious actors may embed malware or sensitive data into DNS records [140]. This becomes even harder to detect when the tunnel is then encrypted, as is the case with DNS over HTTPS [141].

Since its inception, DNS over HTTPS has been adopted by mainstream browsers like Firefox and Chrome to prevent man-in-the-middle attacks and snooping on DNS traffic. While enhancing user privacy, this has contraversely become a field of opportunity for cyber attackers as DNS traffic that is now encrypted may bypass traditional DNS security systems. Thus, the DNS-over-HTTPS (DoH) pipeline can become both a covert channel for malware entering systems as well as data being exfiltrated out. For example, in 2020 the Iranian threat actor group APT34 used DNSExfiltrator2, an open-source DoH tool, to laterally move data to their networks [142]. Malware such as Godlua [143] and PsiXBot [144] have also been spotted in the wild along DoH channels. As an example in our scenario, imagine Bob wants to track activity or exfiltrate data from a server in a competitor’s private network. In order to record this information, he needs to plant a trojan virus on this server. Using DNS-over-HTTPS as a covert data channel, Bob hosts the trojan on a website and uses social engineering to engage a user to click a link, querying the domain and returning the data which contains the trojan. This malware remains undetected as port 443 traffic is not scanned.

7.3.1.4 Covert Data Channels Out

Many enterprise network environments restrict outbound traffic to ports 80 and 443. Users may re-reroute DNS from the standard port 53 to avoid defeat by firewalls, ISPs, or government or enterprise networks [138]. Similarly, firewalls can also be bypassed using SSH over port 443. SSH, or secure shell protocol, gives users a method of accessing another computer over an insecure network. While typically transmitted over port 22, many firewalls default to defeating traffic over this port. Instead, they allow traffic over port 443. Packages exist such as `stunnel` which make SSH tunneling over 443 simple for any Internet user [145]. Creating a covert data channel over SSH can also be a gateway for cache attacks as it is necessary to transmit or exfiltrate the data when accessed [146]. In our scenario, Alice’s middlebox technology is monitoring

traffic on port 22 (SSH) for secure shell sessions. Bob needs to SSH into a machine and copy over some documents. Instead of relying on traditional channels, Bob has re-routed his SSH connection through port 443, creating a covert data channel. Thus, Alice is missing the mission critical data of the file transfer.

7.3.2 System Design

The FORAGER system is a highly adaptable traffic profiling multi-tool capable of representing packets in multi-dimensional space through configurable modules. Users can enable up to four different methods of latent representation for the packets. Each of these modules assesses different portions of the packet data with various hidden representation learning approaches. Specifically, ALPINE and PALM [147] generate one-dimensional locality sensitive hashes which represents single sets of features from the header and payload. MAPLE and DATE generate two-dimensional images and three-dimensional point clouds with clusters, respectively. Votes are provided to the overall framework and used in lookup to generate a classification report. The whole system is depicted in Figure 7.2.

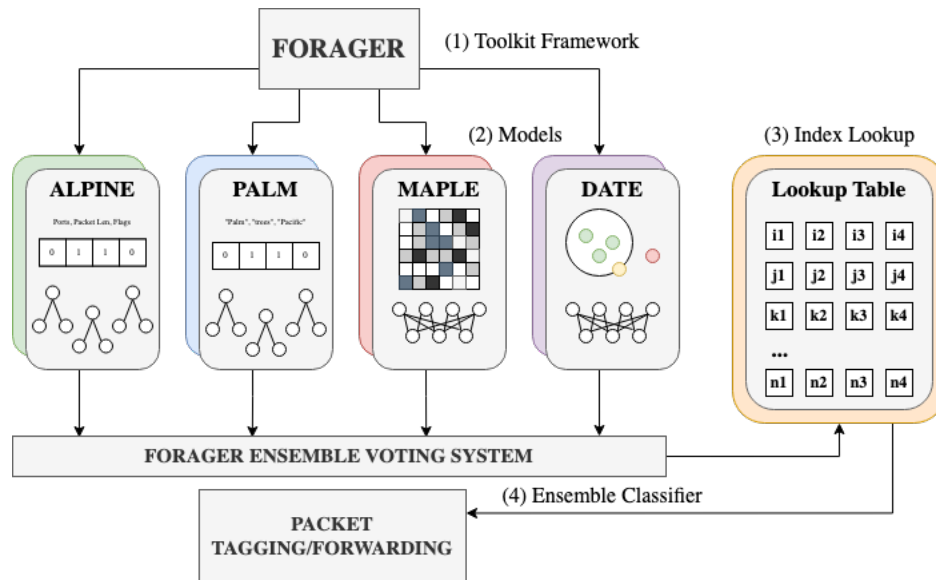


Figure 7.2: System diagram of FORAGER

In the training phase, data is passed in with a label and indexed respective to

class. This index is stored in a lookup table inside FORAGER. Models are enabled or disabled and given a number of votes m per model. This allows the user to more heavily weight one model’s opinion over another, further contributing to the flexibility in real deployments. A few additional hyperparameters may be configured for the neural network models such as epochs for training, and data is processed and weights may be saved for later use. In the testing phase, test data is processed in a similar manner and the input passed through the classifier. DATE amalgamates the votes across models and uses its lookup table to return a classification result. In a real packet processing pipeline, we could add this result as part of metadata or an encapsulation header for the packet for downstream processing. The current implementation generates a report of classification results which may be returned to an administrator for further review and possible identification of traffic on non-standard or unexpected ports which could be malicious.

7.3.3 Datasets

To create a pipeline of data representative of a diverse capture on port 443, we utilized several public datasets created by and widely used in the research community for network traffic classification.

7.3.3.1 ISCX 2016 VPN-nonVPN Dataset

We used the UNB ISCX 2016 VPN-nonVPN dataset [90]. This dataset is a real-world simulation capture created by the Canadian Institute for Cybersecurity (CIC) specifically for encrypted VPN traffic classification [33]. The traffic was captured using Wireshark and tcpdump and generated using OpenVPN; it contains 28GB of traffic. The packets are a mixture of encrypted and compressed or plaintext traffic. Because these are real-world captures, there are some packets such as ICMP and ARP which are service-level and irrelevant for application and traffic type classification. For this reason, we discard non-TCP/UDP packets as part of the process. Previous

work [38, 148, 32] in machine learning-based traffic classification has made use of this public dataset and made similar modifications. A breakdown of the dataset is provided in Table 7.1. The data is divided into the following categories:

- *Traffic Type* - There are seven types of traffic: Browsing, File Transfer, Email, Chat, Streaming, VoIP, and P2P.
- *Application* - Several of the PCAPs were associated specifically with one application such as Skype or Google Hangouts. The applications we uniquely identify are Pidgin (IAM and ICQ chat), BitTorrent (P2P), Facebook, Skype, Google Hangouts, Spotify, Vimeo, Youtube, Netflix, Thunderbird (Email), and Filezilla (FTP).

7.3.3.2 Tor Datasets

The Canadian Institute for Cybersecurity similarly created the UNB ISCX 2016 Tor-nonTor [91] dataset which captures data routed over Tor. The traffic was generated using Whonix ¹ and captured using Wireshark and tcpdump; the Tor dataset contains 22GB of traffic. The benign traffic used for comparison in their traffic profiling work [34] is the same from the VPN dataset, establishing a good baseline for our work in combining these two. A gateway was configured to convert non-encapsulated traffic into Tor PCAPs which were then routed to the Tor network. Because Tor is a circuit-oriented network, all traffic from the gateway to the entry node will be routed through the same connection [34] which could cause over-fitting to a particular dataset. In order to better simulate the diversity of network routing methods, we added as another Tor sample PCAPs from Skynet, a Tor-powered botnet with publicly available datasets for research purposes [89]. For a botnet, the advantage of Tor is hidden services. There is no way to trace the original IP address of a hidden server which is published with its .onion pseudo-domain. All botnet traffic is encrypted.

¹<https://www.whonix.org>

Table 7.1: ISCX VPN/non-VPN and Tor/non-Tor Datasets

Traffic	Content
Browsing	Firefox and Chrome
Email	SMTPS, POP3S and IMAPS
Chat	ICQ, AIM, Skype, Facebook, Hangouts
Streaming	Vimeo, Youtube, Spotify
File Transfer	Skype, FTPS and SFTP using Filezilla
VoIP	Facebook, Skype, Hangouts voice calls
P2P	uTorrent, BitTorrent, Vuze

Application	Benign	VPN	Tor
AIM chat	6K	1.4K	2K
Email	50K	25K	163K
Facebook	1837K	36K	88K
Filezilla (FTPS)	312K	315K	143K
Filezilla (SFTP)	629K	182K	176K
Gmail	15K	–	–
Chrome	–	–	356K
Hangouts	1871K	151K	1114K
Pidgin (ICQ)	4K	7K	2K
Netflix	455K	1433K	–
SCP	777K	–	–
Skype	2261K	1069K	564K
Spotify	63K	193K	118K
Bittorrent	–	133K	–
VoipBuster	1017K	823K	–
Vimeo	214K	591K	288K
Vuze	–	–	264K
Youtube	377K	333K	261K

The malware which infects the host opens a SOCKS proxy on port 55080 reachable through the .onion domain, and can then run a number of bundled attacks and operations like bitcoin mining, data exfiltration, and DDoS attacks [89]. Thus, we leverage evidence of this botnet implementation as an additional Tor data connection.

7.3.3.3 CIRA-CIC-DoHBrw-2020 Dataset

CIC has also made the CIRA-CIC-DoHBrw-2020 dataset containing DNS-over-HTTPS (DoH) and non-DoH traffic on port 443, with a secondary layer of benign versus malicious DoH traffic [140]. The traffic was generated by querying the top 10K

Alexa websites. For non-DoH traffic, an HTTPS request is sent for the website with regular, plaintext DNS. Firefox and Chrome browsers with the configured DNS-over-HTTPS setting were used for benign DoH traffic. For each of the browsers, there were four DNS services used: Cloudflare, AdGuard, Google, and Quad9. Finally, for malicious DoH traffic, tunneling tools dns2tcp, DNSCat2 and Iodine were used to create encrypted TCP covert data channels for transmission. A domain and authoritative name server were established and the DoH tunneling tool used in the client/server interaction. The packet pre-processing we performed on this data was similar to the VPN and Tor dataset processing; a distribution of the data is provided in Table 7.2.

7.3.3.4 EMews SoH Dataset

SSH is another exploitable protocol for both malware injection and data exfiltration over HTTPS. In 2018, Ricks et al [149] created an SSH over HTTPS (SoH) simulation using the eMews framework and CORE to generate packet traces. The network consists of 1022 nodes, 36 of which incorporate the SSH sessions and the others perform various web-crawling activity for benign comparison. EMews ² itself is an open-source, large-scale network traffic data generation tool designed to emulate networks employing protocols which require human interaction (ex: starting an SSH connection). Sample counts for this dataset are provided in Table 7.3.

Table 7.2: Summary of CIRA-CIC-DoHBrw-2020 Dataset

Classification	Application	Sample Count
<i>non-DoH</i>	Chrome	542K
	Firefox	356K
<i>Benign DoH</i>	Chrome	3.5K
	Firefox	16.2K
<i>Malicious DoH</i>	dns2tcp	168K
	DNSCat2	36K
	Iodine	47K

²<https://mews.sv.cmu.edu/research/emews/>

Table 7.3: Summary of eMews SSH over HTTPS Dataset

Classification	Sample Count
HTTPS	542K
SSH over HTTPS	3.5K

7.4 Experiments and Results

We use measures of precision, recall, and F1 score for classifier assessment. In the binary confusion matrix, true positive (TP) indicates correct classification of data as some protocol class C . True negative (TN) is correct classification of data as *not* class C . false negative (FN) implies incorrect identification of traffic as not C , and false positive (FP) is the incorrect classification of data as C .

$$R(C) = \frac{TP}{TP + FP}, P(C) = \frac{TP}{TP + FN}, F1(C) = \frac{2(P \times R)}{P + R} \quad (7.1)$$

A high recall indicates that in the binary problem, much of the traffic is correctly identified as the positive class. High precision in the model shows the quality of positive classifications, i.e. indicates how many were correctly classified out of everything said to be in the positive class. A model with low recall and high precision will likely mis-classify traffic which does belong in the positive class (high rate of false negatives), whereas a model with high recall and low precision will produce more false positives. In the following scenarios, we analyzed different potential traffic classification problems for port 443 based on the described attacks or opportunities from the threat models.

7.4.1 VPN Traffic Profiling

We analyzed the ISCX VPN/non-VPN dataset according to the applications and traffic types identified by the researchers who generated the data and Lotfollahi et al [38] who used the same data for per-packet classification in their Deep Packet models. They used a convolutional neural network and stacked autencoder model for

implementations. We used three different configurations of FORAGER for the traffic and application classification tasks. For each problem, we ran a combination of ALPINE and one of the payload transformation models. For application classification, DeepPacket’s CNN model achieved very high F1 scores for most of the applications. Overall, each model is quite successful in distinguishing applications without inspecting packet contents, preserving user privacy while still accomplishing the mission. In the threat model, Alice could use our tool to better understand what kind of Internet behavior Bob is engaging in over his VPN connection. We notice in both application and traffic classification tasks, our FORAGER models achieve some improved results in the VoIP categories. Our models also perform significantly better in profiling the benign/non-VPN traffic types, up to 0.13 higher than the stacked auto encoder in the VoIP category as well. When ICQ and AIM were combined under the Chat category, FORAGER achieves much higher results than any other model. For all models, distinguishing ICQ and AIM traffic was a difficult task. In further investigation into the network environment, we found both of these traffic streams were generated from the Pidgin messaging application, which makes sense that they were misclassified as one another in all the models. We also note that apps which use similar protocol stacks or likely similar APIs tend to be misclassified as one another. for example, FTPS and SFTP (both using FTP) are more difficult to distinguish from one another than either against SCP using SSH [38]. This is an important note for training and defining classes for neural network models in general, as well-defined classes or clusters make classification a clearer task. One significant difference in the experimental setup of Deep Packet versus FORAGER is that DeepPacket’s CNN requires 300 epochs to train the entire network. MAPLE achieves near comparable results at a mere 10 epochs of training on the same data and split. DATE also achieves good results only training with 10 epochs on its statistical model. We ran additional tests, increasing the number of epochs in order to determine if this improved our models’ efficacy.

In each case, 10 epochs proved to be sufficient for training FORAGER to maximum performance capability.

Table 7.4: Results for classifying VPN data

Application	DeepPacket-CNN			DeepPacket-SAE			ALPINE + PALM			ALPINE + MAPLE			ALPINE + DATE		
	<i>P</i>	<i>R</i>	<i>F1</i>	<i>P</i>	<i>R</i>	<i>F1</i>	<i>P</i>	<i>R</i>	<i>F1</i>	<i>P</i>	<i>R</i>	<i>F1</i>	<i>P</i>	<i>R</i>	<i>F1</i>
AIM chat	0.87	0.76	0.81	0.76	0.64	0.70	0.80	0.55	0.65	0.80	0.57	0.67	0.67	0.69	0.68
Email	0.97	0.82	0.89	0.94	0.99	0.97	0.81	0.76	0.79	0.74	0.80	0.77	0.86	0.73	0.79
Facebook	0.96	0.95	0.96	0.94	0.95	0.95	0.96	0.90	0.93	0.97	0.91	0.94	0.98	0.90	0.94
FTPS	1.00	1.00	1.00	0.97	0.77	0.86	1.00	0.98	0.99	0.99	0.99	0.99	0.99	0.99	0.99
Gmail	0.97	0.95	0.96	0.93	0.94	0.94	0.96	0.87	0.91	0.94	0.90	0.92	0.95	0.90	0.92
Hangouts	0.96	0.98	0.97	0.94	0.99	0.97	0.98	0.89	0.93	0.96	0.91	0.94	0.96	0.91	0.94
ICQ	0.72	0.80	0.76	0.69	0.69	0.69	0.47	0.93	0.63	0.59	0.89	0.71	0.60	0.89	0.72
Netflix	1.00	1.00	1.00	1.00	1.00	1.00	0.99	0.98	0.99	0.99	0.99	0.99	0.99	0.99	0.99
SCP	0.97	0.99	0.98	1.00	1.00	1.00	0.98	0.91	0.94	0.98	0.92	0.95	0.97	0.91	0.94
SFTP	1.00	1.00	1.00	0.70	0.96	0.81	0.99	1.00	1.00	0.99	1.00	0.99	0.99	1.00	1.00
Skype	0.94	0.99	0.97	0.95	0.93	0.94	0.99	0.95	0.97	0.98	0.92	0.95	0.98	0.92	0.95
Spotify	0.98	0.98	0.98	0.98	0.98	0.98	0.99	0.92	0.95	0.99	0.95	0.97	0.99	0.95	0.97
Torrent	1.00	1.00	1.00	1.00	1.00	1.00	0.99	0.97	0.98	0.98	0.99	0.98	0.98	0.99	0.98
VoipBuster	0.99	1.00	0.99	0.99	0.99	0.99	1.00	0.98	0.99	1.00	0.98	0.99	1.00	0.99	0.99
Vimeo	0.99	0.99	0.99	0.99	0.98	0.98	0.99	0.98	0.99	0.97	0.98	0.99	0.98	0.98	0.98
Youtube	0.99	0.99	0.99	0.98	0.99	0.99	0.99	0.94	0.96	0.96	0.96	0.96	0.96	0.96	0.96
Traffic Type	DeepPacket-CNN			DeepPacket-SAE			ALPINE + PALM			ALPINE + MAPLE			ALPINE + DATE		
	<i>P</i>	<i>R</i>	<i>F1</i>	<i>P</i>	<i>R</i>	<i>F1</i>	<i>P</i>	<i>R</i>	<i>F1</i>	<i>P</i>	<i>R</i>	<i>F1</i>	<i>P</i>	<i>R</i>	<i>F1</i>
Chat	0.84	0.71	0.77	0.82	0.68	0.74	0.95	0.81	0.87	0.93	0.64	0.76	0.82	0.72	0.77
Email	0.96	0.87	0.91	0.97	0.93	0.95	0.76	0.99	0.86	0.71	0.74	0.73	0.75	0.67	0.71
File Transfer	0.98	1.00	0.99	0.98	0.99	0.99	0.99	0.91	0.95	0.74	0.93	0.82	0.74	0.93	0.82
Stream	0.92	0.87	0.90	0.82	0.84	0.83	0.99	0.98	0.99	0.99	0.98	0.99	0.99	0.98	0.99
VoIP	0.63	0.88	0.74	0.64	0.90	0.75	0.99	0.95	0.97	0.99	0.94	0.96	0.99	0.95	0.97
VPN: Chat	0.98	0.98	0.98	0.95	0.94	0.94	0.96	0.95	0.96	0.96	0.98	0.97	0.96	0.97	0.97
VPN: Transfer	0.99	0.99	0.99	0.98	0.95	0.97	1.00	0.98	0.99	1.00	0.98	0.99	1.00	0.98	0.99
VPN: Email	0.99	0.98	0.99	0.97	0.93	0.95	1.00	0.98	0.99	0.99	0.98	0.99	0.97	0.99	0.98
VPN: Stream	1.00	1.00	1.00	0.99	0.99	0.99	0.99	0.99	0.99	0.99	0.99	0.99	0.99	0.80	0.89
VPN: Tor-rent	1.00	1.00	1.00	0.99	0.97	0.98	0.99	0.99	0.99	0.98	0.99	0.99	0.83	1.00	0.91
VPN: VoIP	0.99	1.00	1.00	0.99	1.00	0.99	1.00	0.98	0.99	0.99	0.98	0.99	0.99	0.98	0.99

7.4.2 Tor Traffic Profiling

Another type of tunneled traffic we want to identify and profile on port 443 is anything running over Tor. We again ran our three configurations of the FORAGER

Table 7.5: Results for classifying Tor data by application

Application	DeepPacket-CNN			DeepPacket-SAE			ALPINE PALM +			ALPINE MAPLE +			ALPINE DATE +		
	<i>P</i>	<i>R</i>	<i>F1</i>	<i>P</i>	<i>R</i>	<i>F1</i>	<i>P</i>	<i>R</i>	<i>F1</i>	<i>P</i>	<i>R</i>	<i>F1</i>	<i>P</i>	<i>R</i>	<i>F1</i>
AIM chat	–	–	–	–	–	–	0.86	0.53	0.66	0.95	0.50	0.65	0.96	0.49	0.65
Chrome	0.00	0.00	0.00	0.03	0.44	0.06	0.93	0.68	0.78	0.91	0.74	0.81	0.91	0.75	0.82
Email	–	–	–	–	–	–	0.94	0.84	0.89	0.94	0.85	0.89	0.94	0.85	0.90
Facebook	0.10	0.24	0.14	0.06	0.28	0.09	0.97	0.88	0.92	0.97	0.89	0.93	0.97	0.89	0.93
FTPS	–	–	–	–	–	–	0.95	0.99	0.97	0.95	0.99	0.97	0.95	0.99	0.97
Gmail	0.97	0.95	0.96	0.93	0.94	0.94	0.94	0.85	0.89	0.89	0.90	0.89	0.94	0.86	0.90
Hangouts	–	–	–	–	–	–	0.91	0.93	0.95	0.95	0.92	0.93	0.93	0.94	0.94
ICQ	–	–	–	–	–	–	0.40	0.94	0.56	0.51	0.97	0.67	0.48	0.98	0.65
Netflix	–	–	–	–	–	–	0.98	0.98	0.98	0.98	0.99	0.98	0.99	0.98	0.99
SCP	0.97	0.99	0.98	1.00	1.00	1.00	0.99	0.90	0.94	0.96	0.92	0.93	0.96	0.92	0.94
SFTP	–	–	–	–	–	–	1.00	0.94	0.97	0.99	0.94	0.97	0.99	0.94	0.97
Skype	–	–	–	–	–	–	0.97	0.91	0.94	0.95	0.92	0.93	0.96	0.91	0.93
Spotify	–	–	–	–	–	–	0.79	0.91	0.85	0.75	0.95	0.84	0.74	0.95	0.83
VoipBuster	–	–	–	–	–	–	1.00	0.98	0.99	1.00	0.98	0.99	1.00	0.98	0.99
Vimeo	0.44	0.36	0.40	0.05	0.91	0.09	0.98	0.94	0.96	0.95	0.95	0.95	0.97	0.95	0.96
Vuze	–	–	–	–	–	–	0.98	0.72	0.83	0.97	0.73	0.83	0.97	0.71	0.82
Youtube	–	–	–	–	–	–	0.95	0.94	0.94	0.94	0.95	0.95	0.94	0.96	0.95

Table 7.6: Results for classifying Tor data by traffic type

Traffic Type	DIDarknet			J48			ALPINE PALM +			ALPINE MAPLE +			ALPINE DATE +		
	<i>P</i>	<i>R</i>	<i>F1</i>	<i>P</i>	<i>R</i>	<i>F1</i>	<i>P</i>	<i>R</i>	<i>F1</i>	<i>P</i>	<i>R</i>	<i>F1</i>	<i>P</i>	<i>R</i>	<i>F1</i>
Audio	0.92	0.92	0.92	0.97	0.97	0.97	0.99	0.94	0.96	0.99	0.94	0.96	0.99	0.94	0.96
Browsing	0.55	0.47	0.51	0.99	0.99	0.99	0.97	0.94	0.95	0.95	0.97	0.96	0.95	0.97	0.96
Chat	0.90	0.86	0.88	0.93	0.93	0.93	0.92	0.74	0.82	0.84	0.80	0.82	0.93	0.70	0.80
Email	0.66	0.67	0.67	0.93	0.93	0.93	0.90	0.90	0.90	0.91	0.87	0.89	0.92	0.86	0.89
File Transfer	0.74	0.75	0.75	0.98	0.98	0.98	0.99	0.92	0.96	0.99	0.92	0.96	0.99	0.92	0.96
P2P	0.90	0.95	0.93	0.99	0.99	0.99	0.96	0.98	0.97	0.95	1.00	0.97	0.96	0.99	0.97
Video	0.82	0.88	0.85	0.98	0.98	0.98	0.99	0.97	0.98	0.98	0.98	0.98	0.98	0.98	0.98
VOIP	0.58	0.61	0.59	0.99	0.99	0.99	0.72	0.97	0.83	0.82	0.96	0.89	0.71	0.99	0.82

toolkit against the Tor dataset. Deep Packet [38] also ran their SAE and CNN models against this data, but only report a portion of results. We provide a comparison against what they do report for completeness, and note that the FORAGER models are a dramatic improvement in terms of capability. Their model is unable to handle most of the application classes, whereas FORAGER can adapt to all applications and returns results as high as those in the VPN classification task. This result demonstrates the versatility of FORAGER compared to the state-of-the-art and current existing frameworks. Table 7.5 shows results of the application classification task. We also show

comparative results of traffic type classification with work performed by Lashkari et al [7] in their DIDarknet system which uses a two-dimensional convolutional neural network to characterize VPN and Tor traffic. It is important to note that unlike Deep Packet and FORAGER models, DIDarknet’s CNN relies on flow-based features. Both we and Deep Packet prefer a stateless approach for best-effort, per-packet classification and use features of the individual packets themselves. Choorod et al [6] provide results of single packet classification of the Tor dataset using features extracted with Charcount from the Posit toolkit along with a decision tree-based algorithm called J48. This performed remarkably well on the traffic type classification task, and we also compare against their results. They marginally outperform the FORAGER system; but both identified Tor traffic with a much greater accuracy and F1 score than the flow-based system. Choorod et al did not provide results for application-based identification, and it is unclear if their method would generalize to any other problems or has been attempted on any other data sets. The experiment results demonstrate that the system could be used to isolate Tor traffic on port 443, or any traffic stream. In the context of the threat model, Alice can determine if and when Bob is using Tor routing, which is useful for further track and trace. She also has a much greater insight into the traffic type and applications than other models previously gave, advancing the state of the art. A summarization is provided in Table 7.7.

7.4.3 Plaintext Traffic over HTTPS

Our next experiment was to assess the ability of the model to determine if traffic samples on port 443 are truly encrypted packets (TLS) or not. Many firewalls are configured to default-allow port 443 traffic; an attacker may send any content or malware over this port in order to attempt to bypass a firewall. Senders wishing to bypass middlebox technology may send non-TLS data over port 443 to avoid detection. It is thus significant for surveillance and law enforcement operations to monitor traffic on port 443 for non-encrypted data which may be decipherable. IoT

Table 7.7: Results for VPN/Tor classification

Class Set	Accuracy	P	R	F1	κ
Tor	0.99	0.99	0.99	0.99	0.99
VPN	0.97	0.98	0.97	0.97	0.98
Standard	0.99	0.99	0.99	0.99	0.99
Class Set	Accuracy	P	R	F1	κ
VPN/nonVPN	0.97	0.98	0.97	0.97	0.97
Traffic Type	0.86	0.88	0.86	0.86	0.82
Application	0.92	0.93	0.92	0.92	0.82
Class Set	Accuracy	P	R	F1	κ
Tor/nonTor	0.99	0.99	0.99	0.99	0.99
Traffic Type	0.97	0.97	0.97	0.97	0.96
Application	0.89	0.91	0.89	0.89	0.88

devices also can send unencrypted traffic over port 443 [19] which may contain some useful information. Alice finds it incredibly useful to capture all plaintext data where possible as it can contribute to the legal summation of a case. We created two classes of encrypted and non-encrypted traffic, performed random under-sampling, and ran the model, generating results in Figures 7.3, 7.4, and 7.5. The model tends to overselect traffic as encrypted, and underselect traffic which is plaintext. This is not surprising due to the variety of features present in the single plaintext category, versus the lack of common features generally found or uncovered in the encrypted class. In this case, model design may not be well-suited to a binary classification problem because the class *plaintext* is not well-defined. This tuning factor should be considered by users of the FORAGER system during the configuration process.

7.4.4 Compressed Traffic over HTTPS

Compressed versus encrypted traffic both tend to exhibit highly entropic behavior, with the data approaching uniform distribution in more advanced implementations [134]. An attacker may run compressed data over port 443 (or any port) and middleware technology may assume it is encrypted if using statistical tests like Chi-squared or Shannon entropy [150]. Instead, Alice could use our system to isolate

compressed traffic from encrypted traffic, and thus be able to decompress and process that information in follow-on applications. In order to run this experiment, we first ran compression algorithms on the payloads of FTP_DATA packets. We chose this traffic type as the payloads are largely text-based and often sent compressed over the wire. We used three standard compression libraries - `gzip`, `zlib`, and `bz2`. The packets were randomly divided into groups, producing 10,519 sample compressed payloads of each type. The encrypted traffic from the ISCX 2016 VPN-nonVPN dataset [33] was used for the encrypted traffic type. Figure 7.4 shows FORAGER is capable of identifying compressed from encrypted traffic on port 443 or out of any stream (using ALPINE and MAPLE). Furthermore, we can profile down to specific compression type with great accuracy in Figure 7.5. We can also identify plaintext traffic in Figure 7.3. Alice could use this technology to extract compressed data streams and decompress them for further analysis.

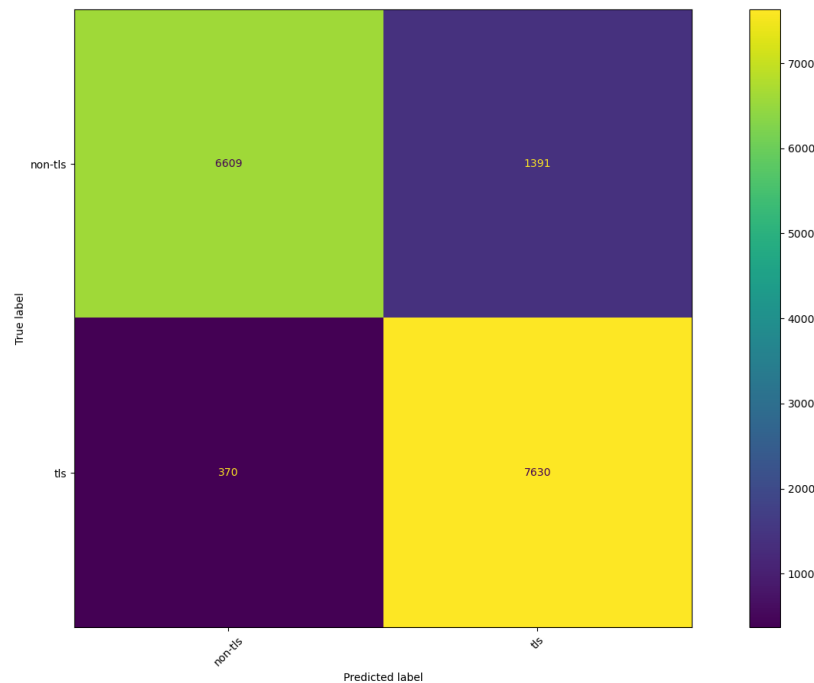


Figure 7.3: Plaintext versus encrypted traffic results

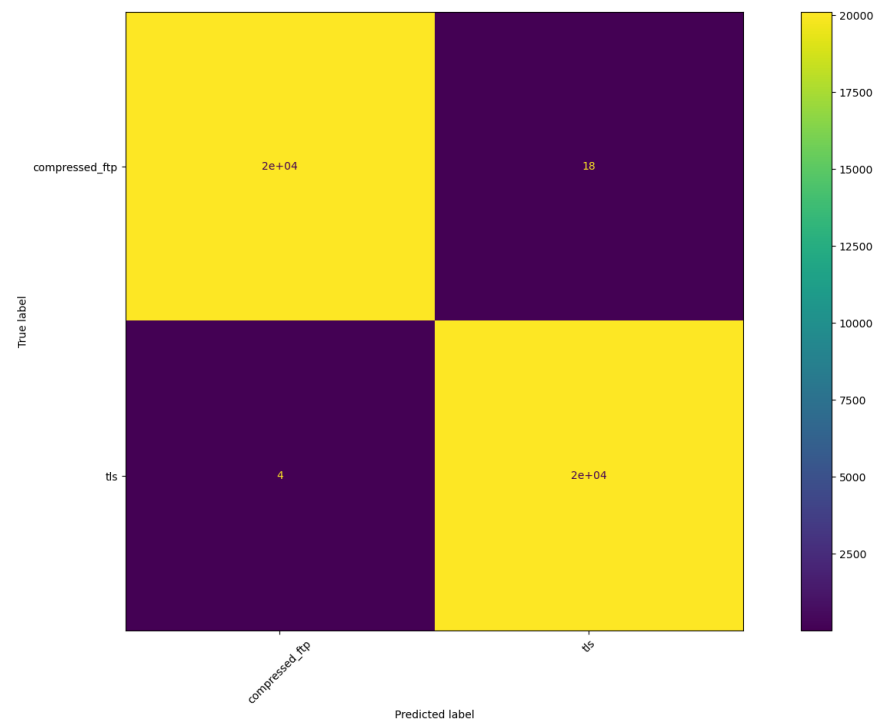


Figure 7.4: Compressed versus encrypted traffic results

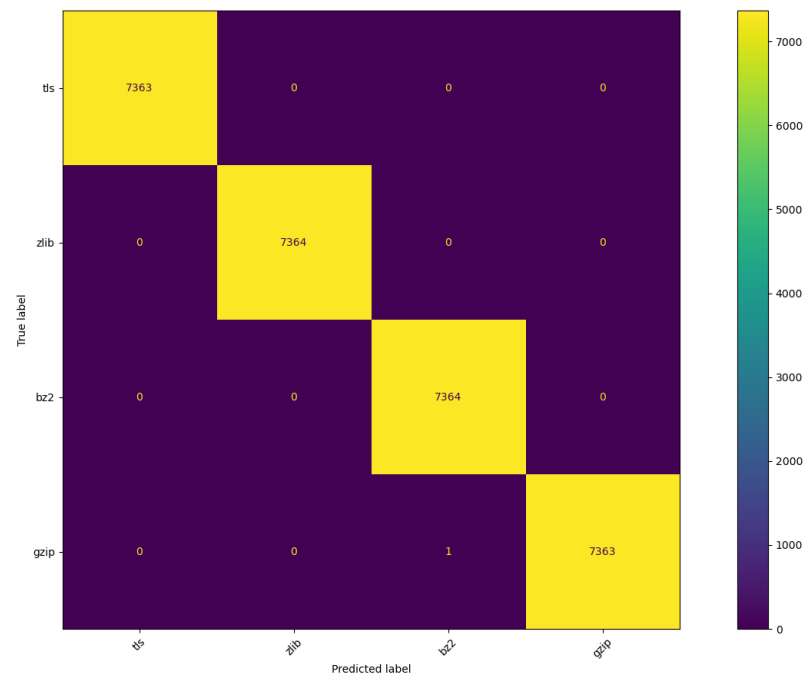


Figure 7.5: Compression type classification results

Table 7.8: Results for encryption and compression detection

Protocol	P	R	F1
Encrypted	1.00	1.00	1.00
Compressed	1.00	1.00	1.00
Encrypted Plaintext	0.85	0.95	0.90
TLS	1.00	1.00	1.00
BZ2	1.00	1.00	1.00
GZIP	1.00	1.00	1.00
ZLIB	1.00	1.00	1.00

7.4.5 SSH over HTTPS

We want to be able to identify particular instances of certain protocols over an encrypted traffic network; One such protocol of interest is SSH. This is often used to bypass firewall restrictions or create covert data channels. In our experiments, we implemented the same model combinations as with the previous scenarios. The FORAGER toolkit was highly successful in classifying SSH traffic from the SSH-over-HTTPS dataset generated by the EMews simulator, achieving over 99% accuracy as shown in Table 7.9. Identifying this traffic would allow inspectors to both mitigate covert data channels or monitor them for possible illicit activity. For example, Alice could detect any open covert data channels Bob might attempt to create over SSH.

7.4.6 DNS over HTTPS

Our final experiment was to detect DNS traffic running over HTTPS. A perpetual challenge for security research in classification and anomaly detection is the notion of intent; just because something is identified as belonging to a particular class or marked as an outlier does not always correlate that it is malicious. For example, a new instance of a running process can just be that - a user running a new program - and not necessarily an executing malware. Or a malformed data packet may appear different than others in the stream, but that could just be a bad transmission. For DNS over HTTPS, mainstream browsers like Chrome and Firefox employ this technique

without malicious intent. As such, being able to determine DNS tunneling as a second layer to the DoH/non-DoH identification is of critical importance. The results in Table 7.9 show that FORAGER rises to meet this challenge, and can distinguish both DoH/non-DoH traffic and then malicious traffic from the identified DoH traffic with high accuracy and few false positives. Thus, any attempt Bob makes at DoH tunneling is now detectable using our toolkit.

Table 7.9: Results for classifying DoH/non-DoH and SoH/non-SoH data

Model	A+P			A+M			A+D		
Class	<i>P</i>	<i>R</i>	<i>F1</i>	<i>P</i>	<i>R</i>	<i>F1</i>	<i>P</i>	<i>R</i>	<i>F1</i>
<i>SSH</i>	1.00	0.99	1.00	0.98	1.00	0.99	1.00	1.00	1.00
<i>Non-SSH</i>	0.99	1.00	1.00	1.00	0.98	0.99	1.00	1.00	1.00
<i>DNS</i>	1.00	0.97	0.98	0.96	0.99	0.98	0.95	0.99	0.97
<i>Non-DNS</i>	0.97	1.00	0.98	0.99	0.96	0.97	0.95	0.97	0.96
<i>Malicious</i>	1.00	0.95	0.97	1.00	0.95	0.97	0.99	0.95	0.97
<i>Benign</i>	0.96	1.00	0.98	0.95	1.00	0.97	0.96	0.99	0.97

7.4.7 Performance Metrics

The experiments were run on a 1.6 GHz Dual-Core Intel Core i5 processor with 16GB of RAM running MacOS Big Sur version 11.6.1. For each of our experiments, we tracked system performance metrics for the models as they performed the classification task and provide this data in Appendix A. This is important for real systems to assess not only the expected tonnage of network traffic they can process, but also to determine what models may be appropriate for deployment. For example, a model whose strategy focuses on online training may not be suitable in an environment who requires high throughput and low downtime, but has a high training time. Or, in some scenarios it may be acceptable to have slightly lower accuracy but a higher throughput rate.

7.5 Discussion and Limitations

Using FORAGER, we are able to provide new insight into a previously under-analyzed gateway into secure networks - port 443. These techniques enable profiling of encrypted traffic streams while preserving user privacy of the content. We can also use FORAGER to identify plaintext or compressed traffic from the encrypted streams which may then be forwarded to decompression and DPI tools. FORAGER can profile the same traffic stream multiple ways, which is an important recognition of the complexity present even in an individual packet. Finally, we show the ability of the system to identify not only particular protocols in the encrypted stream, but also distinguish between benign and malicious instances of their utilization. This would be useful, for example, applied in network intrusion detection (NIDS) and prevention systems. For surveillance missions, we show FORAGER keeps up with the ever-changing, fast-paced real and complex network environment and allows analysts to stay ahead of mission communications and critical data. In seeing through encryption and tunneling, FORAGER proves to be a versatile solution for real network problems.

For plaintext analysis and keyword searching, using regular expression scanning can still be an effective technique. Automatic generation techniques like REXACTOR may also reduce manual overhead required for signature creation and maintenance, and provide useful insight to commonalities across traffic types and protocols. Additional regular expression scanning and generation tools would be a useful addition to the FORAGER toolkit.

CHAPTER 8: CONTRIBUTIONS AND DELIVERABLES

We designed the FORAGER toolkit as a usable software package for reproducibility and further community research. This software is a command line interface (CLI) application installed through PyPi with user documentation provided through ReadTheDocs. The source code repositories are also available through GitHub with open-source licensing. The following components are included in the initial 1.0 release:

TaPCAP: This option parses input PCAP and PCAPNG files, extracting header features and/or payloads into CSV format. This module can be installed separately ¹ or run from FORAGER.

REXACTOR: The token option finds frequent tokens according to a provided threshold from a CSV output from TAPCAP. The regular expression mode performs genetic sequencing alignment of payloads from a CSV output from TAPCAP. REXACTOR can be installed as a separate component ² or run from FORAGER.

ALPINE: In training mode, this model accepts as input the CSV obtained from TAPCAP, generates locality sensitive hashes from the header data, and finalizes and serializes these to an LSHForest [2]. In testing mode, the trained LSHForest is reloaded from the internal cache. Locality sensitive hashes of the input data are generated and index lookup performed. The returned result is applied to the ensemble voting system if other models are configured.

PALM: In training mode, this model accepts as input the CSV obtained from TAPCAP, generates locality sensitive hashes from the payload data, and finalizes and serializes these to an LSHForest [2]. In testing mode, the trained LSHForest

¹<https://pypi.org/project/tapcap/>

²<https://pypi.org/project/rexactor/>

is reloaded from the internal cache. Locality sensitive hashes of the input data is generated and index lookup performed. The returned result is applied to the ensemble voting system if other models are configured.

MAPLE: In training mode, this model accepts as input the CSV obtained from TAPCAP, generates grayscale images from the payload data, and saves the generated neural network model and trained weights to JSON and H5 files, respectively. In testing mode, the trained neural network is reloaded from the internal cache. Grayscale images of the input data are generated and used as input to the CNN classifier. The returned result is then applied to the ensemble voting system if other models are configured.

DATE: In training mode, this model accepts as input the CSV obtained from TAPCAP, generates point clouds from the payload data and performs DBSCAN analysis. It then saves the generated neural network model and trained weights to JSON and H5 files, respectively. In testing mode, the trained neural network and weights are reloaded from the internal cache. Point clouds are generated from input data and DBSCAN analysis performed, and this is used as input to the multi-layer perceptron classifier. The returned result is applied to the ensemble voting system if other models are configured.

Multi-modality: FORAGER supports training and using one or more of the ALPINE, PALM, MAPLE, and DATE models together so that the packet may be analyzed from multiple representations. We use votes in a combined ensemble classifier and return the class with the most votes as the final classification.

CHAPTER 9: CONCLUSIONS

9.1 Future Work

The scope of this dissertation includes the research and methods leading to the initial release of the Forager toolkit. This work focuses on per packet inspection and data transformation techniques which are viable for large-scale deployments. There are a number of flow-based techniques referenced in this work which do have practical applications. Netflow data is often used in recurrent neural network or long short-term memory network analysis and can be extracted often with ease from routers, switches, or other network devices. A future expansion of this toolkit would be the ability to process netflow or other such stream-based data or cross-packet features when available.

There are also many deep learning and data mining techniques which we have not yet covered but have been explored in other works. Principle component analysis (PCA) [151, 152] and eigenvectors [153] have been used in per packet classification with success. These techniques are notably efficient and therefore would be excellent candidates for additional features in a future release for processing at large scale.

While we found that adjusting for TF-IDF score in PALM did not improve classification accuracy in the experiments we ran, there are applications for other natural language processing and word embedding techniques in the literature on traffic classification. Word embeddings can be used to classify real versus fake domain names [154], phishing emails [155, 156, 120] and other kinds of phishing attacks [157], and HTTP traffic [158]. In the cybersecurity domain, Word2Vec has been used to identify cross-site scripting [159, 160] and SQL injection attacks [161] across network packets. Textual analysis is rendered ineffective by encryption and compression in many cases,

but for plaintext analysis or malware detection could be a future approach.

The data itself we chose to work with in these experiments until FORAGER’s initial release has largely been sourced from traditional network devices. Mobile traffic and data gathered from IoT or medical/wearable devices could present unique challenges and opportunities for classification problems. In future versions of this toolkit, we could leverage the dimensionality reduction offered by our data engineering and mining techniques in order to combine feature sets and further increase multimodality. For example, we could embed both the processing power and CPU utilization of an IoT device input along with network traffic embeddings to further inform a classification of IoT devices. There is also much opportunity for experimentation with classifying streams of traffic in multiple ways. As mentioned in the introduction of this work, traffic classification covers a broad scope of problems; importantly, these classes are not distinct from one another and traffic may belong to multiple classes and it may be necessary to classify something multiple times. For example, we may want to determine the user who generated a traffic stream (Bob), what protocols are in the traffic (SIP, RTP), what application Bob used to generate the traffic (Zoom), and what device of Bob’s he used to make the call (MacBook Air). Determining methods to perform multiple classifications like this without re-processing the data each time would be a useful additional contribution.

9.2 Summarization

This work proposes a new system of data mining and deep learning for network traffic classification at scale. The FORAGER system includes models which contribute both novel application of data transformation techniques which unveil hidden representations as well as machine learning models applied in real systems. We propose a method for extracting data from PCAPs and converting them to tabular format, then transforming them appropriately according to what we have learned from previous experiments. For example, plaintext, strongly signed data which we are

aiming to identify by protocol (SIP, HTTP, POP) may be aptly suited for regular expression generation algorithms. Alternatively, weakly signed data may be run with the MAPLE model for payload analysis and ALPINE model for header features for identification. Similarly for profiling, encrypted data may also be well-suited for the ALPINE and MAPLE model combination. Our ensemble voting system allows for this multimodal analysis and can finally be saved and returned through a tagging system implemented in the original data. Thus, our toolkit will be a useful contribution to network analysts and cybersecurity and surveillance specialists for uniquely identifying traffic on the Internet today.

REFERENCES

- [1] F. Yu, Z. Chen, Y. Diao, T. V. Lakshman, and R. H. Katz, “Fast and memory-efficient regular expression matching for deep packet inspection,” in *Proceedings of the 2006 ACM/IEEE Symposium on Architecture for Networking and Communications Systems*, ANCS ’06, (New York, NY, USA), pp. 93–102, Association for Computing Machinery, 2006.
- [2] M. Bawa, T. Condie, and P. Ganesan, “Lsh forest: self-tuning indexes for similarity search,” in *In WWW*, 2005.
- [3] W. Jo, S. Kim, C. Lee, and T. Shon, “Packet preprocessing in cnn-based network intrusion detection system,” *Electronics*, vol. 9, no. 7, p. 1151, 2020.
- [4] T. Project, “Tor project,” 2022.
- [5] E. Jardine, A. M. Lindner, and G. Owenson, “The potential harms of the tor anonymity network cluster disproportionately in free countries,” *Proceedings of the National Academy of Sciences of the United States of America*, vol. 117, pp. 31716 – 31721, 2020.
- [6] P. Choorod and G. Weir, “Tor traffic classification based on encrypted payload characteristics,” in *2021 National Computing Colleges Conference (NCCC)*, pp. 1–6, 2021.
- [7] A. Habibi Lashkari, G. Kaur, and A. Rahali, “Didarknet: A contemporary approach to detect and characterize the darknet traffic using deep image learning,” in *2020 the 10th International Conference on Communication and Network Security*, ICCNS 2020, (New York, NY, USA), pp. 1–13, Association for Computing Machinery, 2021.
- [8] IANA, “Assigned protocol numbers,” February 2021.
- [9] B. Charyyev and M. H. Gunes, “Iot traffic flow identification using locality sensitive hashes,” in *ICC 2020-2020 IEEE International Conference on Communications (ICC)*, pp. 1–6, IEEE, 2020.
- [10] E. P. Freire, A. Ziviani, and R. M. Salles, “Detecting skype flows in web traffic,” in *NOMS 2008 - 2008 IEEE Network Operations and Management Symposium*, pp. 89–96, IEEE, 2008.
- [11] T. Karagiannis, A. Broido, M. Faloutsos, and K. claffy, “Transport layer identification of p2p traffic,” in *Proceedings of the 4th ACM SIGCOMM Conference on Internet Measurement*, IMC ’04, (New York, NY, USA), pp. 121–134, Association for Computing Machinery, 2004.
- [12] M. Smith and R. Hunt, “Network security using nat and napt,” in *Proceedings 10th IEEE International Conference on Networks (ICON 2002). Towards Network Superiority (Cat. No.02EX588)*, pp. 355–360, 2002.

- [13] L. Costa, “Open systems interconnect (osi) model,” *JALA: Journal of the Association for Laboratory Automation*, vol. 3, no. 1, pp. 28–35, 1998.
- [14] J. Zhao, X. Jing, Z. Yan, and W. Pedrycz, “Network traffic classification for data fusion: A survey,” *Information Fusion*, vol. 72, pp. 22–47, 2021.
- [15] M. Hasanzadeh, H. Hamidi, and H. Asghari, “Plaintext transmission over session initiation protocol,” in *7th International Symposium on Telecommunications (IST’2014)*, pp. 629–634, 2014.
- [16] K. A. Ogudo, “Analyzing generic routing encapsulation (gre) and ip security (ipsec) tunneling protocols for secured communication over public networks,” in *2019 International Conference on Advances in Big Data, Computing and Data Communication Systems (icABCD)*, pp. 1–9, 2019.
- [17] S. Zhang, A. Li, H. Zhu, Q. Sun, M. Wang, and Y. Zhang, “Research on the protocols of vpn,” in *Advances in Intelligent Systems and Interactive Applications* (F. Khafa, S. Patnaik, and A. Y. Zomaya, eds.), (Cham), pp. 554–559, Springer International Publishing, 2018.
- [18] D. Maimon, M. Becker, S. Patil, and J. Katz, “Self-protective behaviors over public wifi networks,” in *The LASER Workshop: Learning from Authoritative Security Experiment Results (LASER 2017)*, pp. 69–76, USENIX Association, October 2017.
- [19] D. Wood, N. Apthorpe, and N. Feamster, “Cleartext data transmissions in consumer iot medical devices,” in *Proceedings of the 2017 Workshop on Internet of Things Security and Privacy*, pp. 7–12, 2017.
- [20] M. Capellupo, J. Liranzo, M. Z. A. Bhuiyan, T. Hayajneh, and G. Wang, “Security and attack vector analysis of iot devices,” in *Security, Privacy, and Anonymity in Computation, Communication, and Storage* (G. Wang, M. Atiquzzaman, Z. Yan, and K.-K. R. Choo, eds.), (Cham), pp. 593–606, Springer International Publishing, 2017.
- [21] Y. Wang, E. Kjerstad, and B. Belisario, “A dynamic analysis security testing infrastructure for internet of things,” in *2020 Sixth International Conference on Mobile And Secure Services (MobiSecServ)*, pp. 1–6, 2020.
- [22] M. El, E. McMahon, S. Samtani, M. Patton, and H. Chen, “Benchmarking vulnerability scanners: An experiment on scada devices and scientific instruments,” in *2017 IEEE International Conference on Intelligence and Security Informatics (ISI)*, pp. 83–88, 2017.
- [23] L. Bernaille and R. Teixeira, “Early recognition of encrypted applications,” in *Passive and Active Network Measurement* (S. Uhlig, K. Papagiannaki, and O. Bonaventure, eds.), (Berlin, Heidelberg), pp. 165–175, Springer Berlin Heidelberg, 2007.

- [24] E. Papadogiannaki and S. Ioannidis, “Acceleration of intrusion detection in encrypted network traffic using heterogeneous hardware,” *Sensors*, vol. 21, no. 4, 2021.
- [25] A. Moore, D. Zuev, and M. Crogan, “Discriminators for use in flow-based classification,” 2013.
- [26] M. Roughan, S. Sen, O. Spatscheck, and N. Duffield, “Class-of-service mapping for qos: A statistical signature-based approach to ip traffic classification,” in *Proceedings of the 4th ACM SIGCOMM Conference on Internet Measurement, IMC '04*, (New York, NY, USA), pp. 135–148, Association for Computing Machinery, 2004.
- [27] V. Paxson, “Empirically derived analytic models of wide-area tcp connections,” *IEEE/ACM Transactions on Networking*, vol. 2, no. 4, pp. 316–336, 1994.
- [28] O. Salman, I. H. Eljhajj, A. Kayssi, and A. Chehab, “A review on machine learning-based approaches for internet traffic classification,” *Ann. Telecommun.*, vol. 75, pp. 673–710, December 2020.
- [29] J. Cao, Z. Fang, G. Qu, H. Sun, and D. Zhang, “An accurate traffic classification model based on support vector machines,” *Netw.*, vol. 27, p. n/a, January 2017.
- [30] H.-K. Lim, J.-B. Kim, J.-S. Heo, K. Kim, Y.-G. Hong, and Y.-H. Han, “Packet-based network traffic classification using deep learning,” in *2019 International Conference on Artificial Intelligence in Information and Communication (ICAIIIC)*, pp. 046–051, 2019.
- [31] Z. Li, R. Yuan, and X. Guan, “Accurate classification of the internet traffic based on the svm method,” in *2007 IEEE International Conference on Communications*, pp. 1373–1378, 2007.
- [32] M. Song, J. Ran, and S. Li, “Encrypted traffic classification based on text convolution neural networks,” in *2019 IEEE 7th International Conference on Computer Science and Network Technology (ICCSNT)*, pp. 432–436, 2019.
- [33] G. Draper-Gil., A. H. Lashkari., M. S. I. Mamun., and A. A. Ghorbani., “Characterization of encrypted and vpn traffic using time-related features,” in *Proceedings of the 2nd International Conference on Information Systems Security and Privacy - ICISSP*, pp. 407–414, INSTICC, SciTePress, 2016.
- [34] A. H. Lashkari, G. Draper-Gil, M. S. I. Mamun, and A. A. Ghorbani, “Characterization of tor traffic using time based features,” in *ICISSP*, 2017.
- [35] S. Cui, J. Liu, C. Dong, Z. Lu, and D. DU, “Only header: A reliable encrypted traffic classification framework without privacy risk,” 2022.

- [36] Z. Zou, J. Ge, H. Zheng, Y. Wu, C. Han, and Z. Yao, "Encrypted traffic classification with a convolutional long short-term memory neural network," in *2018 IEEE 20th International Conference on High Performance Computing and Communications; IEEE 16th International Conference on Smart City; IEEE 4th International Conference on Data Science and Systems (HPCC/SmartCity/DSS)*, pp. 329–334, 2018.
- [37] P. Perera, Y.-C. Tian, C. Fidge, and W. Kelly, "A comparison of supervised machine learning algorithms for classification of communications network traffic," in *International Conference on Neural Information Processing*, pp. 445–454, Springer, 2017.
- [38] M. Lotfollahi, R. S. H. Zade, M. J. Siavoshani, and M. Saberian, "Deep packet: A novel approach for encrypted traffic classification using deep learning," *CoRR*, vol. abs/1709.02656, 2017.
- [39] A. Panchenko, L. Niessen, A. Zinnen, and T. Engel, "Website fingerprinting in onion routing based anonymization networks," in *Proceedings of the 10th annual ACM workshop on Privacy in the electronic society*, pp. 103–114, 2011.
- [40] P. Amaral, J. Dinis, P. Pinto, L. Bernardo, J. Tavares, and H. S. Mamede, "Machine learning in software defined networks: Data collection and traffic classification," in *2016 IEEE 24th International conference on network protocols (ICNP)*, pp. 1–5, IEEE, 2016.
- [41] Z. Cao, G. Xiong, Y. Zhao, Z. Li, and L. Guo, "A survey on encrypted traffic classification," in *International Conference on Applications and Techniques in Information Security*, pp. 73–81, Springer, 2014.
- [42] H. A. H. Ibrahim, O. R. A. Al Zuobi, M. A. Al-Namari, G. MohamedAli, and A. A. A. Abdalla, "Internet traffic classification using machine learning approach: Datasets validation issues," in *2016 Conference of Basic Sciences and Engineering Studies (SGCAC)*, pp. 158–166, IEEE, 2016.
- [43] Z. Fan and R. Liu, "Investigation of machine learning based network traffic classification," in *2017 International Symposium on Wireless Communication Systems (ISWCS)*, pp. 1–6, IEEE, 2017.
- [44] G. Sun, T. Chen, Y. Su, and C. Li, "Internet traffic classification based on incremental support vector machines," *Mobile Networks and Applications*, vol. 23, no. 4, pp. 789–796, 2018.
- [45] K.-S. Shim, S.-H. Yoon, S.-K. Lee, and M. Kim, "Sigbox: Automatic signature generation method for fine-grained traffic identification," *J. Inf. Sci. Eng.*, vol. 33, pp. 537–569, 2017.
- [46] K. L. Dias, M. A. Pongelupe, W. M. Caminhas, and L. de Errico, "An innovative approach for real-time network traffic classification," *Computer Networks*, vol. 158, pp. 143–157, 2019.

- [47] Z. Yuan and C. Wang, “An improved network traffic classification algorithm based on hadoop decision tree,” in *2016 IEEE International Conference of On-line Analysis and Computing Science (ICOACS)*, pp. 53–56, 2016.
- [48] E. L. Goodman, C. Zimmerman, and C. Hudson, “Packet2vec: Utilizing word2vec for feature extraction in packet data,” *CoRR*, vol. abs/2004.14477, 2020.
- [49] Y. Wang, Y. Xiang, and S.-Z. Yu, “Automatic application signature construction from unknown traffic,” in *2010 24th IEEE International Conference on Advanced Information Networking and Applications*, pp. 1115–1120, 2010.
- [50] R. Fielding, “Hypertext transfer protocol – http/1.1,” rfc, IETF, June 1999.
- [51] J. R. J. Postel, “File transfer protocol (ftp),” rfc, IETF, October 1985.
- [52] P. Saint-Andre, “Extensible messaging and presence protocol (xmpp): Core,” rfc, IETF, March 2011.
- [53] H.-A. Kim and B. Karp, “Autograph: Toward automated, distributed worm signature detection,” in *Proceedings of the 13th Conference on USENIX Security Symposium - Volume 13*, SSYM’04, (USA), p. 19, USENIX Association, 2004.
- [54] Y. Tang, B. Xiao, and X. Lu, “Using a bioinformatics approach to generate accurate exploit-based signatures for polymorphic worms,” *Computers & Security*, vol. 28, no. 8, pp. 827–842, 2009.
- [55] B. D. Sija, K.-S. Shim, and M.-S. Kim, “Automatic payload signature generation for accurate identification of internet applications and application services,” *KSII Transactions on Internet and Information Systems*, vol. 12, no. 4, pp. 1572–1593, 2018.
- [56] Y. Wang, Y. Xiang, W. Zhou, and S. Yu, “Generating regular expression signatures for network traffic classification in trusted network management,” *Journal of Network and Computer Applications*, vol. 35, no. 3, pp. 992–1000, 2012. Special Issue on Trusted Computing and Communications.
- [57] G. Szabô, Z. Turányi, L. Toka, S. Molnár, and A. Santos, “Automatic protocol signature generation framework for deep packet inspection,” *ICST*, 6 2012.
- [58] B.-C. Park, Y. J. Won, M.-S. Kim, and J. W. Hong, “Towards automated application signature generation for traffic identification,” in *NOMS 2008 - 2008 IEEE Network Operations and Management Symposium*, pp. 160–167, 2008.
- [59] M. Ye, K. Xu, J. Wu, and H. Po, “Autosig-automatically generating signatures for applications,” in *2009 Ninth IEEE International Conference on Computer and Information Technology*, vol. 2, pp. 104–109, 2009.

- [60] C. VinothGeorge and V. Edwards, “Efficient regular expression signature generation for network traffic classification,” 2013.
- [61] H. Toivonen, *Apriori Algorithm*, pp. 60–60. Boston, MA: Springer US, 2017.
- [62] R. Agrawal and R. Srikant, “Fast algorithms for mining association rules in large databases,” in *Proceedings of the 20th International Conference on Very Large Data Bases*, VLDB ’94, (San Francisco, CA, USA), pp. 487–499, Morgan Kaufmann Publishers Inc., 1994.
- [63] S. B. Needleman and C. D. Wunsch, “A general method applicable to the search for similarities in the amino acid sequence of two proteins,” *Journal of Molecular Biology*, vol. 48, no. 3, pp. 443–453, 1970.
- [64] S. Altschul, *basic local alignment search tool (BLAST)*. online, 2003.
- [65] Intel Software, “Hyperscan.”
- [66] X. Wang, Y. Hong, H. Chang, K. Park, G. Langdale, J. Hu, and H. Zhu, “Hyperscan: A fast multi-pattern regex matcher for modern cpus,” in *16th USENIX Symposium on Networked Systems Design and Implementation (NSDI 19)*, (Boston, MA), pp. 631–648, USENIX Association, Feb. 2019.
- [67] tommyod, “efficient-apriori.”
- [68] KimiNewt, “Pyshark.”
- [69] NumPy, “Numpy.”
- [70] Rene Rivera, Bernan Dawes, David Abrahams, “Boost.”
- [71] M. Kapoor, G. Fuchs, and J. Quance, “Rexactor: Automatic regular expression signature generation for stateless packet inspection,” in *2021 IEEE 20th International Symposium on Network Computing and Applications (NCA)*, pp. 1–9, 2021.
- [72] H. Schulzrinne, S. L. Casner, R. Frederick, and V. Jacobson, “RTP: A Transport Protocol for Real-Time Applications.” RFC 3550, july 2003.
- [73] J. Hypolite, J. Sonchack, S. Hershkop, N. Dautenhahn, A. DeHon, and J. M. Smith, *DeepMatch: Practical Deep Packet Inspection in the Data Plane Using Network Processors*, pp. 336–350. New York, NY, USA: Association for Computing Machinery, 2020.
- [74] Google, “Google transparency report,” 2022.
- [75] M. Kapoor, S. Krishnan, and T. Moyer, “Deep packet inspection at scale: Search optimization through locality-sensitive hashing,” in *2022 IEEE 21st International Symposium on Network Computing and Applications (NCA)*, vol. 21, pp. 97–105, 2022.

- [76] Z. Fu, Z. Liu, and J. Li, "Efficient parallelization of regular expression matching for deep inspection," in *2017 26th International Conference on Computer Communication and Networks (ICCCN)*, pp. 1–9, 2017.
- [77] O. Jafari, P. Maurya, P. Nagarkar, K. M. Islam, and C. Crushev, "A survey on locality sensitive hashing algorithms and their applications," 2021.
- [78] P. Indyk and R. Motwani, "Approximate nearest neighbors: Towards removing the curse of dimensionality," in *Proceedings of the Thirtieth Annual ACM Symposium on Theory of Computing, STOC '98*, (New York, NY, USA), pp. 604–613, Association for Computing Machinery, 1998.
- [79] A. Z. Broder, "Identifying and filtering near-duplicate documents," in *Proceedings of the 11th Annual Symposium on Combinatorial Pattern Matching, COM '00*, (Berlin, Heidelberg), pp. 1–10, Springer-Verlag, 2000.
- [80] M. Gabryel, K. Grzanek, and Y. Hayashi, "Browser fingerprint coding methods increasing the effectiveness of user identification in the web traffic," *Journal of Artificial Intelligence and Soft Computing Research*, vol. 10, no. 4, pp. 243–253, 2020.
- [81] U. Bayer, P. M. Comparetti, C. Hlauschek, C. Kruegel, and E. Kirda, "Scalable, behavior-based malware clustering," in *NDSS*, vol. 9, pp. 8–11, Citeseer, 2009.
- [82] P. Lee, L. V. Lakshmanan, and J. X. Yu, "On top-k structural similarity search," in *2012 IEEE 28th International Conference on Data Engineering*, pp. 774–785, 2012.
- [83] H. Yang, Q. He, Z. Liu, and Q. Zhang, "On top-k structural similarity search," *Security and Communication Networks*, vol. 2021, 2021.
- [84] Z. Tang, Q. Wang, W. Li, H. Bao, F. Liu, and W. Wang, "'hslf: Http header sequence based lsh fingerprints for application traffic classification'," in *Paszynski M., Kranzlmüller D., Krzhizhanovskaya V.V., Dongarra J.J., Sloat P.M.A. (eds) Computational Science - ICCS 2021*, 2021.
- [85] W. Jiang and M. Gokhale, "Real-time classification of multimedia traffic using fpga," in *2010 International Conference on Field Programmable Logic and Applications*, pp. 56–63, 2010.
- [86] Y. Leskovec, A. Rajaraman, and J. Ullman, *Mining of Massive Datasets*. Stanford University, 2022.
- [87] J. V. V. Silva, N. R. Oliveira, D. S. V. Medeiros, M. A. Lopez, and D. M. F. Mattos, "A statistical analysis of intrinsic bias of network security datasets for training machine learning mechanisms," February 2022.

- [88] B. Sangster, T. J. O'Connor, T. Cook, R. Fanelli, E. Dean, W. J. Adams, C. Morrell, and G. Conti, "Toward instrumenting network warfare competitions to generate labeled datasets," in *Proceedings of the 2nd Conference on Cyber Security Experimentation and Test*, CSET'09, (USA), p. 9, USENIX Association, 2009.
- [89] C. Guarnieri, "Skynet: a tor-powered botnet straight from reddit," 2012.
- [90] G. D. Gil, A. H. Lashkari, M. Mamun, and A. A. Ghorbani, "VPN/non-VPN Dataset (ISCXVPN2016); Tor/non-Tor Dataset (ISCXTor2016)," 2016.
- [91] G. D. Gil, A. H. Lashkari, M. Mamun, and A. A. Ghorbani, "Tor/non-Tor Dataset (ISCXTor2016)," 2016.
- [92] K. Tumer and J. Ghosh, "Error correlation and error reduction in ensemble classifiers," *Connection science*, vol. 8, no. 3-4, pp. 385–404, 1996.
- [93] TeamStage, "Zoom statistics 2022: How video conferencing changed our lives," 2022.
- [94] D. Curry, "Microsoft teams revenue and usage statistics (2022)," June 2022.
- [95] G. W. Edwards, M. J. Gonzales, and M. A. Sullivan, "Robocalling: Stirred and shaken! - an investigation of calling displays on trust and answer rates," in *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, CHI '20, (New York, NY, USA), pp. 1–12, Association for Computing Machinery, 2020.
- [96] Z. Z. Wu, J. Guo, C. Zhang, and C. Li, "Steganography and steganalysis in voice over ip: A review," *Sensors (Basel, Switzerland)*, vol. 21, 2021.
- [97] S. Nagaraja and R. Shah, "Voiploc: Voip call provenance using acoustic side-channels," *IEEE Security and Privacy 2020*, 2019.
- [98] C. Choti, N. Hnoohom, S. Tritilanunt, and S. Yuenyong, "Prediction of intrusion detection in voice over internet protocol system using machine learning," in *The 2021 9th International Conference on Computer and Communications Management*, pp. 149–155, 2021.
- [99] M. Kmet, P. Matousek, and O. R. Martin, "Fast rtp detection and codecs classification in internet traffic," 2014.
- [100] M. Sha, T. Manesh, and S. M. A. El-atty, "Voip forensic analyzer," *International Journal of Advanced Computer Science and Applications*, vol. 7, 2016.
- [101] D.-Y. Kao, T.-C. Wang, and F.-C. Tsai, "Forensic artifacts of network traffic on wechat calls," in *2020 22nd International Conference on Advanced Communication Technology (ICACT)*, pp. 262–267, IEEE, 2020.

- [102] J. Gu, Z. Wang, J. Kuen, L. Ma, A. Shahroudy, B. Shuai, T. Liu, X. Wang, G. Wang, J. Cai, *et al.*, “Recent advances in convolutional neural networks,” *Pattern recognition*, vol. 77, pp. 354–377, 2018.
- [103] K. O’Shea and R. Nash, “An introduction to convolutional neural networks,” *ArXiv*, vol. abs/1511.08458, 2015.
- [104] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [105] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” *CoRR*, vol. abs/1512.03385, 2015.
- [106] H.-K. Lim, J.-B. Kim, J.-S. Heo, K. Kim, Y.-G. Hong, and Y.-H. Han, “Packet-based network traffic classification using deep learning,” in *2019 International Conference on Artificial Intelligence in Information and Communication (ICAIIIC)*, pp. 046–051, 2019.
- [107] L. Yang and A. Shami, “A Transfer Learning and Optimized CNN Based Intrusion Detection System for Internet of Vehicles,” *arXiv e-prints*, p. arXiv:2201.11812, Jan. 2022.
- [108] J. Gao, “Network intrusion detection method combining cnn and bilstm in cloud computing environment,” vol. 2022, p. 11, 2022.
- [109] C. Tu, E. Takeuchi, A. Carballo, and K. Takeda, “Point cloud compression for 3d lidar sensor using recurrent neural network with residual blocks,” in *2019 International Conference on Robotics and Automation (ICRA)*, pp. 3274–3280, 2019.
- [110] J.-L. Costeux, F. Guyard, and A.-M. Bustos, “Qrp08-5: Detection and comparison of rtp and skype traffic and performance,” in *IEEE Globecom 2006*, pp. 1–5, 2006.
- [111] N. Patwari, A. O. Hero III, and A. Pacholski, “Manifold learning visualization of network traffic data,” in *Proceedings of the 2005 ACM SIGCOMM workshop on Mining network data*, pp. 191–196, 2005.
- [112] Y. Wang and J. M. Solomon, “Deep closest point: Learning representations for point cloud registration,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
- [113] M. Quach, G. Valenzise, and F. Dufaux, “Folding-based compression of point cloud attributes,” in *2020 IEEE International Conference on Image Processing (ICIP)*, pp. 3309–3313, 2020.

- [114] E. Schubert, J. Sander, M. Ester, H. P. Kriegel, and X. Xu, “Dbscan revisited, revisited: why and how you should (still) use dbscan,” *ACM Transactions on Database Systems (TODS)*, vol. 42, no. 3, pp. 1–21, 2017.
- [115] T. Mullin, “Dbscan parameter estimation using python,” july 2020.
- [116] N. Rahmah and I. S. Sitanggang, “Determination of optimal epsilon (eps) value on dbscan algorithm to clustering data on peatland hotspots in sumatra,” in *IOP conference series: earth and environmental science*, vol. 31, p. 012012, IOP Publishing, 2016.
- [117] H. Song and J. W. Lockwood, “Efficient packet classification for network intrusion detection using fpga,” in *FPGA ’05*, 2005.
- [118] C. Wang, M. Ji, J. Wang, W. Wen, T. Li, and Y. Sun, “An improved dbscan method for lidar data segmentation with automatic eps estimation,” *Sensors*, vol. 19, no. 1, 2019.
- [119] T. G. Dietterich, “Ensemble methods in machine learning,” in *International workshop on multiple classifier systems*, pp. 1–15, Springer, 2000.
- [120] F. Tajaddodianfar, J. W. Stokes, and A. Gururajan, “Texception: A character/word-level deep learning model for phishing url detection,” in *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 2857–2861, 2020.
- [121] Y. Lee, H. Kwon, S.-H. Choi, S.-H. Lim, S. H. Baek, and K.-W. Park, “Instruction2vec: Efficient preprocessor of assembly code to detect software weakness with cnn,” *Applied Sciences*, vol. 9, no. 19, 2019.
- [122] H. L. Duarte-Garcia, C. D. Morales-Medina, A. Hernandez-Suarez, G. Sanchez-Perez, K. Toscano-Medina, H. Perez-Meana, and V. Sanchez, “A semi-supervised learning methodology for malware categorization using weighted word embeddings,” in *2019 IEEE European Symposium on Security and Privacy Workshops (EuroS&PW)*, pp. 238–246, IEEE, 2019.
- [123] G. Ayoade, K. A. Akbar, P. Sahoo, Y. Gao, A. Agarwal, K. Jee, L. Khan, and A. Singhal, “Evolving advanced persistent threat detection using provenance graph and metric learning,” in *2020 IEEE Conference on Communications and Network Security (CNS)*, pp. 1–9, IEEE, 2020.
- [124] F. Palau, C. A. Catania, J. L. Guerra, S. Garcia, and M. Rigaki, “DNS tunneling: A deep learning based lexicographical detection approach,” *CoRR*, vol. abs/2006.06122, 2020.
- [125] K. Kishioka, K. Hongyo, T. Kimura, T. Kudo, Y. Inoue, and K. Hirata, “Prediction method of infection spreading with cnn for self-evolving botnets,” in *2018 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, pp. 1810–1815, 2018.

- [126] O. Boyaci, M. R. Narimani, K. R. Davis, M. Ismail, T. J. Overbye, and E. Serpedin, "Joint detection and localization of stealth false data injection attacks in smart grids using graph neural networks," *IEEE Transactions on Smart Grid*, vol. 13, no. 1, pp. 807–819, 2022.
- [127] K. Nagaraj, A. Starke, and J. McNair, "Glass: a graph learning approach for software defined network based smart grid ddos security," in *ICC 2021-IEEE International Conference on Communications*, pp. 1–6, IEEE, 2021.
- [128] L. A. Iliadis and T. Kaifas, "Darknet traffic classification using machine learning techniques," in *2021 10th International Conference on Modern Circuits and Systems Technologies (MOCASST)*, pp. 1–4, 2021.
- [129] H. Ma, J. Cao, B. Mi, D. Huang, Y. Liu, and Z. Zhang, "Dark web traffic detection method based on deep learning," in *2021 IEEE 10th Data Driven Control and Learning Systems Conference (DDCLS)*, pp. 842–847, 2021.
- [130] D. Sarkar, P. Vinod, and S. Y. Yerima, "Detection of tor traffic using deep learning," in *2020 IEEE/ACS 17th International Conference on Computer Systems and Applications (AICCSA)*, (Los Alamitos, CA, USA), pp. 1–8, IEEE Computer Society, nov 2020.
- [131] G. R. S. Weir, "The posit text profiling toolset," 2007.
- [132] W. M. Shbair, T. Cholez, J. FranÃ§ois, and I. Chrisment, "A survey of https traffic and services identification approaches," *ArXiv*, vol. abs/2008.08339, 2020.
- [133] A. P. Singh and M. Singh, "A comparative review of malware analysis and detection in https traffic," *International Journal of Computing and Digital Systems*, vol. 10, no. 1, pp. 111–123, 2021.
- [134] F. Casino, K.-K. R. Choo, and C. Patsakis, "Hedge: Efficient traffic classification of encrypted and compressed packets," *IEEE Transactions on Information Forensics and Security*, vol. 14, no. 11, pp. 2916–2926, 2019.
- [135] D. Hahn, N. J. Aphthorpe, and N. Feamster, "Detecting compressed cleartext traffic from consumer internet of things devices," *ArXiv*, vol. abs/1805.02722, 2018.
- [136] M. Zain ul Abideen, S. Saleem, M. Ejaz, and R. Soundrapandiyan, "Vpn traffic detection in ssl-protected channel," *Sec. and Commun. Netw.*, vol. 2019, jan 2019.
- [137] V. Ghi  tte, H. J. Griffioen, and C. Doerr, "Fingerprinting tooling used for ssh compromisation attempts," in *International Symposium on Recent Advances in Intrusion Detection*, 2019.

- [138] S. Kerr, L. Song, and R. Wan, “A review of DNS over port 80/443,” Internet-Draft draft-shane-review-dns-over-http-04, Internet Engineering Task Force, November 2016. Work in Progress.
- [139] M. Simioni, P. Gladyshev, B. Habibnia, and P. R. Nunes de Souza, “Monitoring an anonymity network: Toward the deanonymization of hidden services,” *Forensic Science International: Digital Investigation*, vol. 38, p. 301135, 2021.
- [140] M. MontazeriShatoori, L. Davidson, G. Kaur, and A. Habibi Lashkari, “Detection of doh tunnels using time-series classification of encrypted traffic,” in *2020 IEEE Intl Conf on Dependable, Autonomic and Secure Computing, Intl Conf on Pervasive Intelligence and Computing, Intl Conf on Cloud and Big Data Computing, Intl Conf on Cyber Science and Technology Congress (DASC/PiCom/CBDCCom/CyberSciTech)*, pp. 63–70, 2020.
- [141] P. E. Hoffman and P. McManus, “DNS Queries over HTTPS (DoH).” RFC 8484, Oct. 2018.
- [142] QuoIntelligence, “How dns-over-https (doh) has changed the threat landscape for companies,” February 2021.
- [143] TrendMicro, “New godlua backdoor found abusing dns over https (doh) protocol,” July 2019.
- [144] P. T. I. Team, September 2019.
- [145] W. Wong, “Stunnel: Sslng internet services easily,” *SANS Institute*, November, 2001.
- [146] C. Maurice, M. Weber, M. Schwarz, L. Giner, D. Gruss, C. A. Boano, S. Mangard, and K. Römer, “Hello from the other side: Ssh over robust cache covert channels in the cloud,” in *Network and Distributed System Security Symposium*, 2017.
- [147] M. Kapoor, S. Krishnan, and T. Moyer, “Deep packet inspection at scale: Search optimization through locality-sensitive hashing,” in *2022 IEEE 21st International Symposium on Network Computing and Applications (NCA)*, vol. 21, pp. 97–105, 2022.
- [148] K. Zhou, W. Wang, C. Wu, and T. Hu, “Practical evaluation of encrypted traffic classification based on a combined method of entropy estimation and neural networks,” *Etri Journal*, vol. 42, pp. 311–323, 2020.
- [149] B. Ricks, P. Tague, and B. Thuraishingham, “Ddos-as-a-smokescreen: Leveraging netflow concurrency and segmentation for faster detection,” in *2021 Third IEEE International Conference on Trust, Privacy and Security in Intelligent Systems and Applications (TPS-ISA)*, pp. 217–224, 2021.

- [150] F. D. Gaspari, D. Hitaj, G. Pagnotta, L. D. Carli, and L. V. Mancini, “Encod: Distinguishing compressed and encrypted file fragments,” in *International Conference on Network and System Security*, pp. 42–62, Springer, 2020.
- [151] R. Yan and R. Liu, “Principal component analysis based network traffic classification,” *Journal of Computers*, vol. 9, pp. 1234–1241, May 2014.
- [152] N. Jadav, N. Dutta, H. K. D. Sarma, E. Pricop, and S. Tanwar, “A machine learning approach to classify network traffic,” in *2021 13th International Conference on Electronics, Computers and Artificial Intelligence (ECAI)*, pp. 1–6, 2021.
- [153] Q. Li, Y. Ju, C. Zhao, and X. He, “An encrypted field locating algorithm for private protocol data based on data reconstruction and moment eigenvector,” *IEEE Access*, vol. 9, pp. 42947–42958, 2021.
- [154] R. Vinayakumar, M. Alazab, S. Srinivasan, Q.-V. Pham, S. K. Padannayil, and K. Simran, “A visualized botnet detection system based deep learning for the internet of things networks of smart cities,” *IEEE Transactions on Industry Applications*, vol. 56, no. 4, pp. 4436–4456, 2020.
- [155] R. S. Rao, A. Umarekar, and A. R. Pais, “Application of word embedding and machine learning in detecting phishing websites,” *Telecommunication Systems*, vol. 79, no. 1, pp. 33–45, 2022.
- [156] H. Yuan, Z. Yang, X. Chen, Y. Li, and W. Liu, “Url2vec: Url modeling with character embeddings for fast and accurate phishing website detection,” in *2018 IEEE Intl Conf on Parallel & Distributed Processing with Applications, Ubiquitous Computing & Communications, Big Data & Cloud Computing, Social Computing & Networking, Sustainable Computing & Communications (ISPA/IUCC/BDCloud/SocialCom/SustainCom)*, pp. 265–272, IEEE, 2018.
- [157] M. Somesha and A. R. Pais, “Classification of phishing email using word embedding and machine learning techniques,” *Journal of Cyber Security and Mobility*, pp. 279–320, 2022.
- [158] M. Gniewkowski, H. Maciejewski, T. R. Surmacz, and W. Walentynowicz, “Http2vec: Embedding of http requests for detection of anomalous traffic,” *arXiv preprint arXiv:2108.01763*, 2021.
- [159] Y. Fang, Y. Li, L. Liu, and C. Huang, “Deepxss: Cross site scripting detection based on deep learning,” in *Proceedings of the 2018 international conference on computing and artificial intelligence*, pp. 47–51, 2018.
- [160] Z. Guichang, N. Yunfeng, and W. Xing, “Cnnpayl: An intrusion detection system of cross-site script detection,” in *Proceedings of the 2019 11th International Conference on Machine Learning and Computing*, pp. 477–483, 2019.

- [161] L. Yu, S. Luo, and L. Pan, “Detecting sql injection attacks based on text analysis,” in *3rd International Conference on Computer Engineering, Information Science & Application Technology (ICCIA 2019)*, pp. 95–101, Atlantis Press, 2019.