



THE UNIVERSITY *of* EDINBURGH

This thesis has been submitted in fulfilment of the requirements for a postgraduate degree (e.g. PhD, MPhil, DClinPsychol) at the University of Edinburgh. Please note the following terms and conditions of use:

This work is protected by copyright and other intellectual property rights, which are retained by the thesis author, unless otherwise stated.

A copy can be downloaded for personal non-commercial research or study, without prior permission or charge.

This thesis cannot be reproduced or quoted extensively from without first obtaining permission in writing from the author.

The content must not be changed in any way or sold commercially in any format or medium without the formal permission of the author.

When referring to this work, full bibliographic details including the author, title, awarding institution and date of the thesis must be given.

Traffic microstructures and network anomaly detection

Henry Clausen

A thesis presented for the degree of
Doctor of Philosophy



THE UNIVERSITY of EDINBURGH
informatics

Laboratory for Foundations of Computer Science (LFCS)

School of Informatics

University of Edinburgh

March 1, 2022

Abstract

Much hope has been put in the modelling of network traffic with machine learning methods to detect previously unseen attacks. Many methods rely on features on a microscopic level such as packet sizes or interarrival times to identify reoccurring patterns and detect deviations from them. However, the success of these methods depends both on the quality of corresponding training and evaluation data as well as the understanding of the structures that methods learn. Currently, the academic community is lacking both, with widely used synthetic datasets facing serious problems and the disconnect between methods and data being named the “semantic gap”.

This thesis provides extensive examinations of the necessary requirements on traffic generation and microscopic traffic structures to enable the effective training and improvement of anomaly detection models. We first present and examine *DetGen*, a container-based traffic generation paradigm that enables precise control and ground truth information over factors that shape traffic microstructures. The goal of DetGen is to provide researchers with extensive ground truth information and enable the generation of customisable datasets that provide realistic structural diversity.

DetGen was designed according to four specific traffic requirements that dataset generation needs to fulfil to enable machine-learning models to learn accurate and generalisable traffic representations. Current network intrusion datasets fail to meet these requirements, which we believe is one of the reasons for the lacking success of anomaly-based detection methods. We demonstrate the significance of these requirements experimentally by examining how model performance decreases when these requirements are not met.

We then focus on the control and information over traffic microstructures that DetGen provides, and the corresponding benefits when examining and improving model failures for overall model development. We use three metrics to demonstrate that DetGen is able to provide more control and isolation over the generated traffic. The ground truth information DetGen provides enables us to probe two state-of-the-art traffic clas-

sifiers for failures on certain traffic structures, and the corresponding fixes in the model design almost halve the number of misclassifications .

Drawing on these results, we propose *CBAM*, an anomaly detection model that detects network access attacks through deviations from reoccurring flow sequence patterns. *CBAM* is inspired by the design of self-supervised language models, and improves the AUC of current state-of-the-art by up to 140%. By understanding why several flow sequence structures present difficulties to our model, we make targeted design decisions that improve on these difficulties and ultimately boost the performance of our model.

Lastly, we examine how the control and adversarial perturbation of traffic microstructures can be used by an attacker to evade detection. We show that in a stepping-stone attack, an attacker can evade every current detection model by mimicking the patterns observed in streaming services.

Lay Summary

Computers in a network communicate by exchanging packets that each contain a piece of data. Different applications and settings use different mechanisms to package this information, which results in persistent structures in the overall packet stream.

Attackers abuse flaws in the mechanisms of applications that authenticate users, and can thus access and manipulate information that they should not have access to. Most of these attacks are executed remotely by sending packets with malicious code to the victim. Signature detection systems, the predominant defence against these attacks, recognise small pieces of data in the packet stream that are characteristic for a known attack. Albeit very efficient, this method does not protect from novel attacks for which this information is not available yet.

Attacks that are inserted and executed in a packet stream often introduce perturbations to the persistent structures that are observed in the packet streams of benign communication. Network anomaly detection aims to train machine learning methods on benign traffic streams that are capable of recognising attacks as deviations from the learned packet stream structures in order to detect attacks without requiring any information about the attack. A main difficulty so far has been the quality of available data to examine the nature of these structures, which has resulted in network anomaly detection not yet producing detection rates sufficient for operational deployment.

This thesis examines how these traffic structures are formed, how they can be generated in a controlled way, and how an improved understanding traffic structures can boost the results and design of network anomaly detection methods. For this, we present and examine DetGen, a traffic generation paradigm that provides more control than current methods to enable precise examinations of how detection models react to specific traffic structures.

We also present CBAM, a novel network anomaly detection method that identifies dependencies in a traffic stream in a similar way as we identify dependencies between words in a sentence. By examining and understanding which dependencies are harder

to learn for the model, we are able to improve the model design and ultimately outperform the state-of-the art in the detection of seven network attacks.

Lastly, we examine how well an attacker with control can evade existing detection methods when adding altering traffic patterns with adversarial noise. It turns out that for a specific attack, an attacker can successfully evade all existing methods.

Declaration

I declare that this thesis was composed by myself, that the work contained herein is my own except where explicitly stated otherwise in the text, and that this work has not been submitted for any other degree or professional qualification. Some of the material used in this thesis has been published in the following papers:

- Henry Clausen, Robert Flood, and David Aspinall: “Traffic generation using containerization for machine learning”. In: *Proceedings of the 2019 ACSAC DYNAMICS Workshop* (appears in Chapter 3) [1].
- Henry Clausen and David Aspinall: “Examining traffic microstructures to improve model development”. In: *2021 IEEE Security and Privacy Workshops* (appears in Chapter 4) [2].
- Henry Clausen, David Aspinall, and Rober Flood: “Controlling network traffic microstructures for machine-learning model probing”. In: *Proceedings of the 2021 EAI SecureComm Conference* (appears in Chapter 4) [3].
- Henry Clausen, Gudmund Grov, Marc Sabate, and David Aspinall: “Better Anomaly Detection for Access Attacks Using Deep Bidirectional LSTMs”. In: *Machine Learning for Networking: Third International Conference (MLN), 2020* (appears in Chapter 5) [4].
- Henry Clausen, Gudmund Grov, and David Aspinall: “CBAM: A Contextual Model for Network Anomaly Detection”. In: *Computers 10.6 (2021) (MPDI)* (appears in Chapter 5) [5].
- Henry Clausen, Michael S. Gibson, and David Aspinall: “Evading stepping-stone detection with enough chaff”. In: *International Conference on Network and System Security (NSS), 2020* (appears in Chapter 6) [6].

Henry Clausen

Acknowledgements

This thesis would have not been possible without the constant help from my supervisor David Aspinall. We spent countless hours seeking research possibilities, planning projects, and discussing improvements of results. I am especially grateful for his support that reached across countries and for making me better at articulating and communicating my sometimes muddled thoughts.

I also want to thank my industry supervisors Michael Gibson and Alex Healing, whose support was vital throughout this PhD. Their industry perspective helped me focus on problems that have relevant application instead of losing myself in details. I enjoyed collaborating with them on two projects, and I received a glimpse into how cyber-defense research is conducted in industry during my time in Ipswich.

This PhD has been kindly and generously supported in part by *BT group* through their PhD scholarship program, and in part by the *EPSRC* through their industry CASE conversion scheme. I would also like to thank the *School of Informatics* and the *Laboratory for Foundations of Computer Science*, who enabled me to conduct my PhD in an inspirational environment with great facilities, and supported me on various travels.

I also want to thank Robert Flood, whom I enjoyed tremendously to collaborate with on several projects, and who had significant input in the implementation of Det-Gen. I want to thank both him and Alec Fren for providing invaluable input when I needed last minute feedback for paper submissions on too many occasions.

I am very grateful to Gordon Ross, Mahesh Marina, Rik Sarkar, and Paul Patras, who each helped me find a new perspective to the results of this PhD and focus on the branches that produce the most significant contributions.

I would not have been able to complete this chapter of my life without the support of my friends and family. I want to specially thank Leonie for the mutual encouragement as well as distraction during my time in Edinburgh, Benedikt and Marius for great discussions about flaws of publishing in academia. Lastly, I want to thank my parents Katrin and Gerhard as well as my siblings Gloria, Stella, and Victor who always had my back when I struggled.

Table of Contents

1	Introduction	19
1.1	Traffic microstructures	21
1.1.1	Application of traffic microstructure models	22
1.2	Contributions	22
1.3	Central research objectives	24
1.4	Thesis overview	25
2	Background	27
2.1	Network traffic and attacks	27
2.1.1	Traffic metadata	27
2.1.2	Network attacks	28
2.2	Anomaly detection	29
2.2.1	Supervised vs. Unsupervised learning	30
2.2.2	Language models and masked language modelling	31
2.3	Network intrusion detection	33
2.3.1	Deep vs shallow packet inspection	33
2.3.2	Misuse vs anomaly detection	34
2.3.3	Evolution of network anomaly-detection	35
2.4	Traffic datasets	37
2.4.1	Real-world datasets	37
2.4.2	Synthetic datasets	39
2.4.3	Generative traffic models and traffic generation	41
3	Requirements for Machine Learning	45
3.1	Introduction	45
3.1.1	Outline	46
3.1.2	Scope of DetGen and potential use-cases	47

3.2	Background	48
3.2.1	Containerisation with Docker	48
3.3	Problems in modern datasets	51
3.3.1	Dataset Requirements	53
3.4	Design	54
3.4.1	Scenarios	55
3.4.2	Subscenarios	55
3.4.3	Randomisation within Subscenarios	56
3.4.4	Network transmission	56
3.4.5	Capture	57
3.4.6	Implementation Process	57
3.4.7	Simple Example Scenario - FTP server	58
3.4.8	Dataset creation	59
3.4.9	Implemented scenarios	60
3.5	Validation experiments	60
3.5.1	Reproducible scenarios	62
3.5.2	Exploring Artificial Delays	63
3.5.3	Advantages of Dynamic Dataset Generation	66
3.6	Conclusions	69
3.6.1	Difficulties and limitations	70
3.6.2	Future work	70
4	Traffic generation to probe and understand model behaviour	73
4.1	Introduction	73
4.1.1	Existing datasets and ground-truth information	75
4.1.2	Outline	76
4.2	Methodology and example	76
4.3	Refining the notion of benign traffic for anomaly detection	79
4.3.1	Projection coherency evaluation	80
4.3.2	Investigating individual cluster incoherences	82
4.4	Traffic microstructures and their influence factors	84
4.5	How DetGen generates precisely controlled data	86
4.5.1	Simulation of external influence	87
4.6	Verifying the generative determinism of DetGen	88
4.7	Reconstructing an IDS-dataset for efficient probing	91

4.8	Conclusions	93
4.8.1	Difficulties and limitations:	94
5	CBAM: An anomaly detection model for traffic microstructures	95
5.1	Introduction	95
5.1.1	Outline	96
5.2	Overview	97
5.3	Evaluation pitfalls	99
5.4	Design	100
5.4.1	Session construction	100
5.4.2	Contextual modelling	100
5.4.3	Architecture selection	102
5.4.4	Trained architecture	103
5.4.5	Parameter selection and training	106
5.4.6	Detection method	107
5.5	Datasets and benchmark models	107
5.5.1	Dataset assembly	107
5.5.2	Dataset split	110
5.5.3	Sample imbalance and evaluation methodology	111
5.5.4	Benchmark comparison models	112
5.6	Detection performance	113
5.6.1	CICIDS-17 results	113
5.6.2	LANL results	116
5.6.3	How attacks affect flow structures	118
5.6.4	Runtime performance	119
5.7	Benign traffic and longterm stability	120
5.7.1	UGR-16 data	120
5.7.2	CICIDS-17 and LANL-15 results	123
5.7.3	Importance of training data size	123
5.8	Benefit of increased model complexity	125
5.8.1	Bidirectionality for better session context	125
5.8.2	Additional layers for complex session modelling	127
5.8.3	Comparison with simpler models	129
5.9	Related work	129
5.10	Limitations and evasion	131

5.10.1	Limitations	131
5.10.2	Evasion and resilience	132
5.10.3	Future Work	133
5.11	Conclusion	134
6	Stepping-stone detection and evasive microstructure control	135
6.1	Introduction	135
6.1.1	Outline	136
6.1.2	Background and threat model	137
6.2	Related work	138
6.2.1	Testbeds and data	138
6.2.2	Detection methods and threat models	139
6.3	Data generation setting	139
6.3.1	Containerisation with DetGen	139
6.3.2	Simulating stepping stones with SSH-tunnels and Docker	139
6.3.3	Evasive tactics	141
6.4	Evaluation data	142
6.4.1	Stepping-stone data	142
6.4.2	Benign data	142
6.4.3	Evaluation methodology	143
6.5	Selected SSD methods and Implementation	144
6.6	Results	146
6.6.1	Data without evasion tactics	146
6.6.2	Delays	147
6.6.3	Chaff	148
6.6.4	False positives	149
6.6.5	Influence of chain length	150
6.6.6	Influence of network settings	151
6.6.7	Summary	151
6.7	Conclusion	151
7	Conclusion	153
7.1	Revisiting research objectives	153
7.2	Critical analysis	155
7.2.1	Usage of synthetic data	155
7.2.2	Operational considerations	156

7.2.3	Potential incompleteness of microstructure control	156
7.3	Future work	156
7.3.1	Future application of DetGen	156
7.3.2	Model and data adaptivity	157
7.3.3	Further application of NLP-models and large-scale datasets . .	158
7.3.4	Directions for detecting stepping-stones	158
7.4	Final outlook	159
Bibliography		161

Acronyms

NID	network intrusion detection
IDS	intrusion detection system
DPI	deep packet inspection
SPI	shallow packet inspection
ML	machine learning
NLP	natural language processing
LH	likelihood
LLH	log-likelihood
AS	anomaly score
FP	false positive
FPR	false positive rate
TPR	true positive rate
IAT	interarrival time
SSD	stepping-stone detection
RTT	round-trip time
RNN	recurrent neural network
LSTM	long-short-term memory
CNN	convolutional neural network
SVM	support vector machine
GAN	generative adversarial network

Chapter 1

Introduction

Intrusion detection can be seen as a never ending arms-race. On the one side, malicious actors that aim to access, manipulate, or damage computation infrastructure in a network. And on the other detection system providers and research community that provide tools to defend against attacks. An important part of this defence are network intrusion detection (NID) systems, which analyse computer traffic for to detect intrusion attempts. An increasing threat to today's computer networks is the usage of zero-day exploits in attacks that circumvent traditional NID models. A 2019 Cyber Assessment Framework guidance from the UK National Cyber Security Centre (NCSC) states that “*some cyber attackers will go to great lengths to avoid detection via standard security monitoring tools such as anti-virus software, or signature-based network intrusion detection systems, which give a direct indication of compromise*” [7].

Traditionally the most common approach to NID is based on detecting closely defined notions of attacks with what are known as *attack signatures*. Since the late 1980's, signature based systems such as Snort [8] have dominated the field of network intrusion detection due to high effectiveness, low false positives, and good computational efficiency. Due to their design, signature-based methods can only alert on known issues that had been categorised as threats on a signature list; zero-day attacks remain a blind spot of traditional NID systems. With attackers having access to more resources than ever, custom-built malware is becoming more common to circumvent signature-based detection. A report from Watchguard Technologies found that 74% of malware attacks in 2021 were using zero-day exploits [9].

In the 1990's, *network anomaly detection* emerged as a complementary detection tool. Through statistical analysis or the training of machine learning (ML) methods, anomaly detection learn behaviour structures in benign traffic in order to identify attack

behaviour violating these structures as deviations from learned benign behaviour. By not relying on specific signatures of attacks, anomaly methods can detect new attacks and therefore decrease the threat of zero-day attacks. Even though signature detection remains the predominant technique in NID systems, anomaly-based detection solutions are finding more widespread industrial application. Products from companies such as CrowdStrike and Darktrace [10, 11] promise protection from zero-day-threats through anomaly-detection, and the NCSC Cyber Assessment Framework guidance states: “*The science of anomaly detection, which goes beyond using pre-defined or prescriptive pattern matching, is a challenging but growing area. Capabilities like machine learning are increasingly being shown to have applicability and potential in the field of intrusion detection.*” However, the absence of suitable agreed benchmarks makes it difficult to compare products, and the NCSC has recently issued guidance on the use of Intelligent Security Tools after concerns that companies purchase such tools without always understanding their capabilities [12].

In academia, successful research results for anomaly-based NID so far has been restricted to the detection of high-volume attacks. Recent evaluations show that the current anomaly-based network intrusion detection methods fail to detect remote access attacks reliably [13]. Already in 2010, Paxson and Sommer have identified what they called a *semantic gap* between structures in network traffic that are significant for detection, and the design of machine learning models aimed to perform network anomaly detection [14].

In our eyes, one reason for this gap is the fact that no effort has been made so far to monitor or control the various factors that shape traffic structures. The current quasi-benchmark NID-datasets CICIDS-17/18, UNSW-NB-15 and KDD-99 pay more attention to the inclusion of specific attacks and topologies rather than the documentation of the generated traffic. This situation has so far led researchers to often simply evaluate a variety of ML-models on these datasets in the hope of edging out competitors, sometimes without fully understanding model or data flaws.

This thesis introduces the notion of **traffic microstructures** and addresses the *semantic gap* between these structures and corresponding model design and improvement. For this, we set out to control, manipulate, and control important microstructures to monitor their effect on anomaly detection models. The aim of this thesis is to understand traffic microstructures and the factors that shape them, and to define methodologies to develop machine-learning based anomaly detection methods that leverage traffic microstructures effectively for the detection of access attacks.

1.1 Traffic microstructures

Both computational scalability concerns and the increasing encryption of traffic led researchers to shift the application of ML methods from packet payloads to packet and connection metadata. Prominent recent network intrusion detection methods such as Kitsune or DeepCorr learn structures in the sizes, flags, or interarrival times (IATs) of packet sequences for decision-making [15, 16]. These structures reveal information about attack behaviour, but are also influenced by a number of factors such as network congestion or the software version. We define these microscopic structures as follows:

Traffic microstructures are reoccurring patterns in the metadata and temporal ordering of packet sequences in an individual connection, such as the packet sizes of a Diffie-Hellman exchange or typical IATs of video streaming, or across a short sequence of connections on a host, such as the pattern of port 53 (DNS) connections being followed by port 80 or 443 connections (HTTP/S). Traffic microstructures are shaped by and characteristic for the specific communication activity as well as the computational conditions that facilitate the communication.

Learning accurate representations of these microstructures is a key to build successful traffic anomaly models. However, controlling or monitoring traffic microstructures has not been considered so far in the generation of datasets. Labelling traffic samples is labour intensive and often impossible without additional information about the underlying computational processes. This leaves researchers often in the dark about systematic model-failures, or understanding specific properties of traffic that is not processed correctly.

Furthermore, privacy and security concerns so far have prevented the release of large and detailed real-world datasets suitable for network intrusion detection. Researchers are therefore forced to turn to synthetically generated datasets. The current NID traffic generation processes however provides far less variability and realism than necessary for effective training of machine-learning methods.

In comparison to other domains where machine-learning models have been continuously developed and improved, the success of machine-learning based NID methods has not progressed as much, with researchers mostly trying to apply new techniques to the same datasets rather than understanding and improving existing methods. This lack of progress has been described in a qualitative manner by Paxson and Sommer in 2010 and confirmed by several papers since then, and a 2018 survey by Nisioti et al. has given evidence that models detecting access attacks continue to produce high

false positive rates in the area of up to several percent. [14, 17, 18, 13]. We argue that improving the understanding of traffic microstructures and the factors that shape them as well as extensive traffic control, monitoring and labelling is one of the keys to improve the design and performance of anomaly-based NID methods.

1.1.1 Application of traffic microstructure models

The communication between two applications is a result of computational activities that can normally be described as a *finite-state machine*. Models of traffic microstructures should model the packet and connection sequences resulting from this communication under the influence of external influences (transmission failures, computational load, etc.) and the content and complexity of transmitted data. A well-built model should therefore be able to distinguish different computational activities or applications or identify deviations from the communication of known activities or applications, and can therefore be used to classify traffic and identify attacks on network applications. For this, the attacks within the detection scope should take place at or close to the network application layer to generate deviations in the observed traffic microstructures. A privilege escalation attack on a host executed via a remote SSH-shell session would hardly leave a visible deviation in the observed microstructures as the performed SSH-communication and corresponding packet streams still conform to the SSH-protocol, whereas a direct exploit on the SSH-application during a session would likely cause more observable microstructure deviations.

1.2 Contributions

This thesis examines the necessary requirements on **traffic generation** and **traffic microstructures** to enable the effective training and improvement of **anomaly detection models**. Fig. 1.1 highlights these three main themes that this thesis revolves around, and their examined interactions.

We first study how current machine-learning based detection methods are designed and their requirements on comprehensive and realistic datasets to generate models that generalise well on persistent traffic microstructures. We establish our findings with three experiments and empirical evidence to demonstrate how the non-compliance with these requirements affects model performance.

The thesis then proposes DetGen, a traffic generation paradigm to provide near

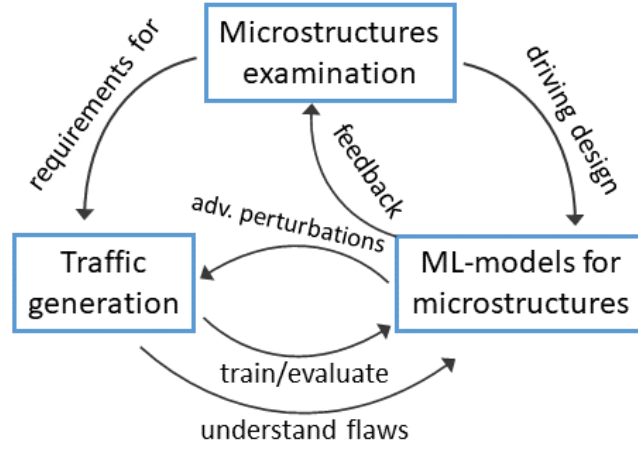


Figure 1.1: Three main themes of this thesis.

deterministic control over factors that shape traffic microstructures. The goal of DetGen is to provide researchers with extensive ground truth information and generate customisable, scalable, and controllable datasets with realistic levels of microstructure diversity. We verify the level of control DetGen provides with two experiments and compare the results with existing traffic generation techniques.

We then demonstrate how the precise traffic control and corresponding ground truth labels that DetGen provides can improve model development with two experimental studies. For this, we probe two state-of-the-art NID models with traffic traces that are subject to controlled variations to understand where the models fail and how they can be improved. We identify design flaws with regards to network retransmissions and HTTP-requests, and are thus able to improve the detection performance to almost halve the number of misclassifications with very simple design adjustments. These studies serve as examples of how the model development process in NID can be expanded to improve existing models in new ways.

Considering these results, we propose CBAM, a **contextual bidirectional anomaly model** that detects low-volume access attacks as deviations from recurring flow microstructures. CBAM is trained as a self-supervised traffic language model and outperforms the detection rates of current state-of-the-art methods by up to 140%. We examine how persistent microstructures are perturbed by specific attacks, and how CBAM detects these perturbations. The output of CBAM also provides valuable insights into the frequency, clustering, and evolution of microstructures in real-world traffic data.

Lastly, we look at stepping-stone attacks where an attacker relays an attack via

a jump host located in the target network. We examine how attackers can perturb the relayed traffic in an adversarial way to prevent models from identifying the corresponding relay of persistent microstructures. We use DetGen to generate controllable stepping-stone attack data to identify which level of perturbations are required to prevent detection. We then evaluate several models with different detection methodologies on this data and conclude that unfortunately realistic adversarial perturbations can evade all current methods. This suggests future directions for necessary research on adversarial robustness that build up on our work.

1.3 Central research objectives

The research contained in this thesis was conducted under the following objectives:

Research Objective 1

How well-structured is the space of microstructures observed in the traffic of a machine or a network? To what degree are these microstructures a result of specific computational activities that are of interest for traffic classification and network intrusion detection, and how much are they affected by other external variables?

Research Objective 2

How can we identify microstructures significant for intrusion detection and train corresponding models? What requirements must a labelled traffic generation framework fulfil to provide realistic data?

Research Objective 3

To what degree can relevant microstructures in network traffic be captured in a model from a training dataset, and how can we achieve this? How can a model adapt to changes of structures in benign traffic?

Research Objective 4

To what degree can access attacks be detected by a model that learns traffic microstructures? What kind of attacks will necessarily show contextual anomalies, and which will not? Can an adversary adapt his attacks to avoid detection?

We will revisit these objectives and how well they were addressed in this thesis in Section 7.1.

1.4 Thesis overview

This section explains the structure that this thesis follows. Fig. 1.2 provides a graphical overview over the thesis themes discussed in each chapter.

Chapter 2 provides the reader with background information on network traffic and attacks, machine-learning based network intrusion detection, and explains the necessary concepts of machine-learning based anomaly detection for sequential data streams. The chapter highlights important NID models that leverage traffic microstructures, and discusses existing datasets and synthetic traffic generation as well as their current shortcomings.

Chapter 3 introduces DetGen, a framework that provides small-scale traffic generation specifically aimed for machine-learning based intrusion detection. The chapter starts with the identification of four major shortcomings of existing NID-datasets regarding the training of machine-learning models that leverage traffic microstructures. The chapter then defines four requirements that a machine-learning-focused traffic generation framework must fulfil. The chapter provides results from several experiments to demonstrate that models do not sufficiently generalise when the training data does not meet these requirements.

Chapter 4 builds up on the results in Chapter 3 by experimentally examining the influence of several factors on traffic microstructures and the level of control that DetGen provides. The chapter then provides two case-studies on how to produce traffic effectively to probe a state-of-the-art traffic classifier, and why a high degree of generative determinism is required for this to isolate the influence of traffic microstructures.

Chapter 5 introduces CBAM, an anomaly detection model that detects low-volume access attacks as deviations from persistent flow microstructures. The chapter explains the reasoning for CBAM's design in order to learn specific microstructures in a self-supervised manner, and how attacks perturb these microstructures. The chapter then provides detection results and a comparison to three benchmark models, before examining long-term stability of the learnt structures and current shortcomings on real-world datasets.

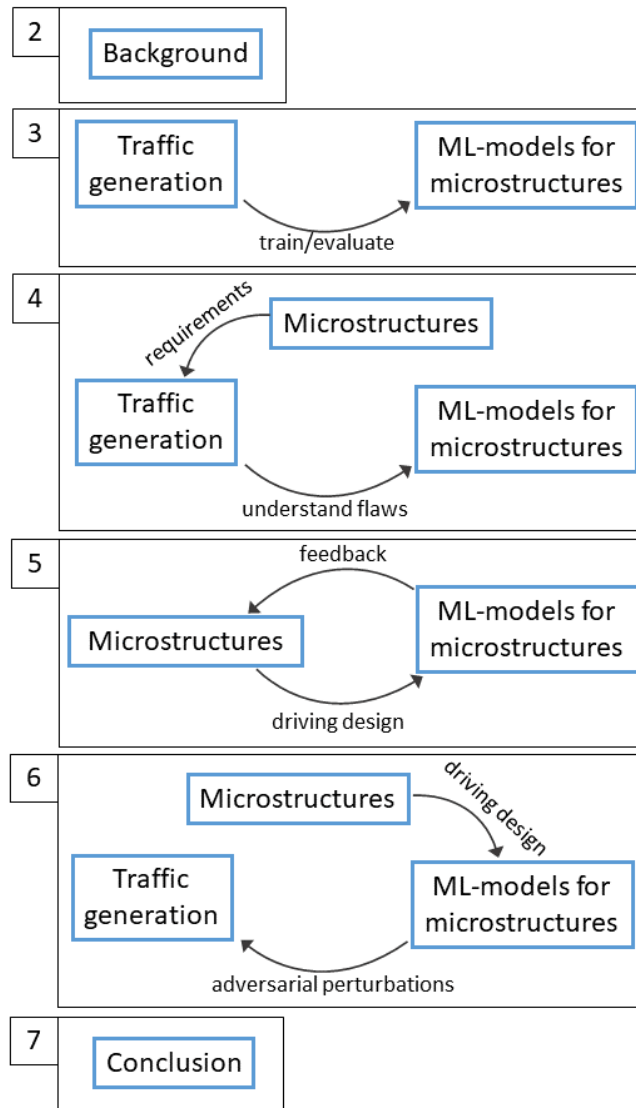


Figure 1.2: Chapter structure of this thesis.

Chapter 6 explains the procedure of a stepping-stone attack, common techniques to detect them, and how attackers introduce microscopic perturbations to evade detection. The chapter then provides results based on data that we generated with varying degrees of introduced perturbations to compare the performance of seven stepping-stone detection methods.

Chapter 7 concludes this thesis, provides a critical perspective on the presented results, and gives an outlook to future work in this area.

Chapter 2

Background

2.1 Network traffic and attacks

2.1.1 Traffic metadata

Computers in a network mostly communicate by sending *network packets* to each other, which contain the control information necessary for transmission, called the packet header, and the user information, called payload. Information contained in the header depends on the respective transport layer protocol (TCP, UDP, ICMP, etc.). Fig. 2.1 displays the content of a TCP-header as an example.

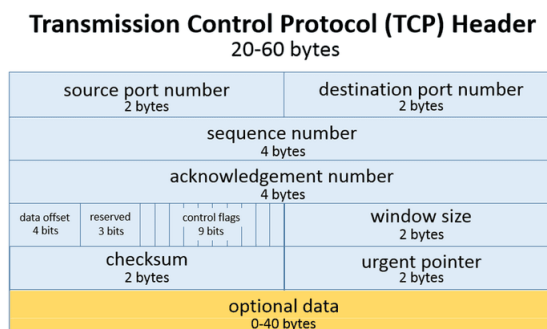


Figure 2.1: Typical format of a TCP packet. Source: <https://www.lifewire.com/tcp-headers-and-udp-headers-explained-817970>, accessed 23.08.2021

In a monitoring setting, packets are usually captured by network routers and stored in the *pcap* format. In case of space shortage or privacy concerns, packet streams are often summarised into **network flows**. RFC 3697 [19] defines a network flow as a sequence of packets that share the same source and destination IP addresses, IP protocol, and for TCP and UDP connections the same source and destination ports. A network

flow is usually saved containing this information along with the start and duration of the connection and the number of packets and bytes transferred in the connection, but can also contain additional information describing the connection. Both data formats are used widely in network intrusion detection.

2.1.2 Network attacks

One reason for the recent rise of cyber crime is the increased use of sophisticated techniques for the attack of specific targets. Attackers use social engineering and custom-build malware that penetrate defensive barriers to stay undetected in an infected system for an extended period of time. In 1980, James P. Anderson, a member of the *Defense Science Board Task Force on Computer Security* at the U.S. Air Force, published the first report to introduce the notion of automated intrusion detection [20]. In it, he defines an **intrusion attempt** as

“...an unauthorized and deliberate attempt to access or manipulate information, or to render a system unreliable or unusable.”

Very often, intrusive attacks involve some sort of network communication between the victim machine(s) and a malicious agent. As this thesis focuses primarily on network intrusions and corresponding defence systems, we will take a closer look at this type of communication. A recent survey covering intrusive attacks and defence systems distinguishes four classes of malicious network attacks [13]:

1. *DoS-attacks*: A denial-of-service attack is an attempt to remove the ability of a particular computer to communicate with other machines over an extended period of time. All major types of DoS-attacks overwhelm the target server with requests that are corrupted to bind server resources.
2. *Network probing attack*: The purpose of network probing is to gather information about computers in a network. This typically involves sending specific service requests to other computers in the network to find out about open ports or the operating system running on a machine.
3. *Access Attacks*: These are attacks that attempt to gain unauthorized access directly to a machine or to resources on a machine. This is typically done by brute-forcing the authentication system or exploiting vulnerabilities to execute code on the targeted machine. This could both be an individual from outside

gaining access to the network, or a user from inside the network accessing services or privileges outside of their authority. Examples of such attacks are SQL-injections, where information is extracted from an SQL-server through the covert execution of SQL-commands in malicious requests, or by planting and executing software scripts on a machine using techniques such as cross-site-scripting or buffer-overflow exploits. This thesis is primarily concerned with the detection of network access attacks.

4. *Data Manipulation Attack*: Also known as “man-in-the-middle”, these attacks involve an attacker reading and manipulating information in a data stream that is not addressed to them by exploiting authentication mechanisms. A common form of such an exploit is *IP spoofing* where an attacker pretends to be a trusted computer by sending packets with a spoofed trusted source IP address.

2.2 Anomaly detection

Anomaly detection refers to the problem of identifying data instances that have significantly different properties than previously observed “normal” data instances. Anomaly detection has its origins in statistical research where normal data instances are assumed to be generated as random variables from a probability distribution $P_N(X)$. A new data sample X_{new} is then identified as anomalous if its properties correspond to regions with vanishing likelihood, i.e. $P_N(X_{\text{new}}) \approx 0$. Usually, the hard part in anomaly detection is to use observed data efficiently to build an estimated distribution $\hat{P}_N(X)$ that resembles $P_N(X)$, and often relies on training machine-learning models.

Anomalies can be classified into a) point-, b) group-, or c) contextual/behavioural anomalies. Point anomalies or outliers occur when a single data sample differs from the rest of the samples in a dataset through its position in the feature space, such as an unusually large or long connection in a network. Group anomalies occur when the occurrence frequency of a feature subset is anomalous within an observation interval, such as the frequency of half-open connections during a DoS attack. Contextual or behavioural anomalies occur when one or more data samples do not conform to expected behaviour within their temporal or spatial context, such as unusually large SQL-connections being opened after an HTTP-request from an unauthorised host during an SQL-injection attack. Graphical examples of point and contextual anomalies are depicted in Fig. 2.2.

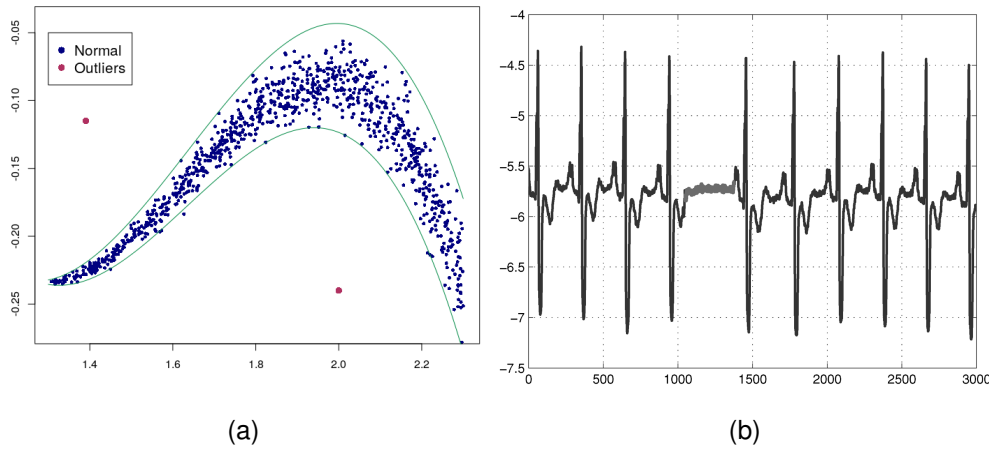


Figure 2.2: The left plot (a) depicts simple anomalies that deviate in distance from regular data. The right plot (b) shows a contextual anomaly where a group of data instances do not follow a repeating timely behaviour with respect to all other data points (corresponding to an *Atrial Premature Contraction*, taken from [21])

2.2.1 Supervised vs. Unsupervised learning

Machine learning methods are often grouped into supervised methods trained on labelled data, and unsupervised methods that are trained on unlabelled data.¹ Besides the input data samples, supervised methods receive separate response data that the model is then trained on to predict in a general settings without any available response data. Typically, the response data is a class label for the corresponding data sample, or numerical values in a regression-type problem. Unsupervised methods are trained entirely on the input samples in the absence of a separate response. The purpose is often not the prediction itself, but to extract information such as numerical thresholds, sample likelihoods, or groupings that describe the data. Examples of unsupervised methods are clustering methods or time-series and sequence prediction.

Machine-learning-based anomaly-detection typically deal with a two-class classification problem, where one class, benign data, is abundant and coherent in their structure while the other, anomalous data samples, are rare, incoherent, and thus more difficult to describe. Anomaly-detection methods are therefore mostly trained in an unsupervised manner entirely on benign data samples to extract information on their underlying structures and provide boundaries of typical behaviour to distinguish them from anomalous samples that exceed these boundaries.

¹Semi-supervised methods are sometimes identified as a separate group, but their functionality is often closely related to supervised methods.

Real-world benign data typically exhibits structures in a myriad of ways, and it is often impossible to build a model that reflects all of these structures in an accurate way. Even though anomaly-detection methods are often claimed to be independent of a notion of anomalous or attack behaviour, the design of the method and the choice of structures that the learned benign-boundary is generated from, mean that some knowledge of anomalous behaviour is needed.

Anomaly detection has found wide application in areas where stability of critical systems or detection of their abuse is crucial. Examples of these include engine-fault detection in wind-turbines and fraud detection for credit cards. The assumption here is that the modelled system, such as the sensor-data stream of an engine or the buying behaviour of a customer, remains stable and generates consistent data.

2.2.2 Language models and masked language modelling

In self-supervised learning, a subgroup of unsupervised approaches, an artificial response is often generated from the input data to generate a labelled dataset from an unlabelled dataset to train a model on large amounts of unlabelled data to learn general representations of data structures before fine-tuning it on smaller labelled dataset for a specific task, such as classifying the topic of a text or producing a caption for an image. A common way to achieve this is by masking some part of the input and training the model to predicting it from the remaining parts of the input, also called self-prediction. Common examples of this approach include the masking of pixel rows which are then predicted with PixelCNN models using an approach called *contrastive predictive coding* [22], or by masking individual words in an input sentence that predicted from the remaining sentence, a technique called *masked language modelling* that is used by self-supervised language models such as ALBERT [23] and that we will examine in more detail now.

In natural language processing (NLP), language models are used for text completion, generation, and correction. A statistical language model is a probability distribution that describes the overall likelihood of a sequences of words. For example, in the sentence “*The students opened their laptops.*”, the language model could provide the probability

$$\Pr (“laptops” | “The students opened their”)$$

of the word “laptop” following the previous part of the sentence.

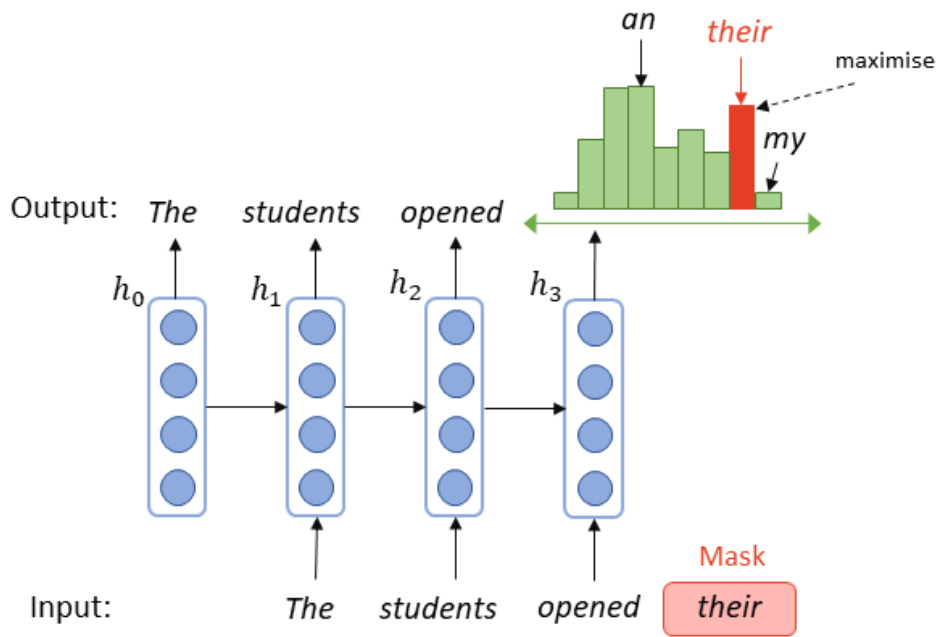


Figure 2.3: Depiction of the self-supervised training of a language model.

Language models are often trained using *masked language modelling* with large amounts of unlabelled text. The model is fed an input sentence with a part of the sentence masked that the model tries to predict. The predictions are then compared to the actual words, as depicted in Fig. 2.3.

In the past, language models were mainly built using simple *N-gram* models or *Hidden Markov* models. In the last decade the success of deep learning pushed *LSTM-based* (Long-short-term memory) networks as depicted in Fig. 2.3 to the top of well-performing language models. Recently, LSTM-networks have been overtaken by *Transformer*-networks that can be trained more efficiently on huge datasets.

Language models can be used as anomaly detectors in for example a grammar checking setting. Here, the predicted probability of a word or a whole sentence are compared to a threshold value, and grammatical or syntactical anomalies are detected as low predicted probabilities.

Written language and network traffic metadata bear significant similarity as each consists of sequences of discrete events (words vs. packets/flows) that stand in contextual relationship to their surroundings. However, the application of successful language models to network anomaly detection is still unexplored. A part of this thesis tries to demonstrate the potential of training language models to capture traffic microstructures.

2.3 Network intrusion detection

The field of intrusion detection is concerned with the development of methods and tools that identify and locate possible intrusions in a computer network. Intrusion detection is a well researched area, with the first IDSs emerging in the late 1980's. Intrusion detection today comprises a variety of research areas in terms of different types of data sources, system architectures, detection scope, and so forth. In 1987, Dorothy Denning established the notion of two different types IDS implementations [24]: a) host-based and b) network-based. *Host-based intrusion detection systems* monitor each host in a network individually, and rely on local data streams such as application logs or raw operating system calls as data source. *Network-based intrusion detection* refers to the detection of malicious traffic in a network of computers and uses network traffic captures as a data source. Host-based systems have the advantage of working with high quality data that are typically very informative. NIDS have the advantage of being platform independent and more robust against attacks on the NIDS as detection of an attack is not done on the attacked system. Fig 2.4 provides an overview over different aspects to consider in network intrusion detection.

2.3.1 Deep vs shallow packet inspection

Modern firewalls and NID solutions rely extensively on *deep packet inspection* (DPI), where the payloads of packets are scanned for attack signatures or anomalies such as known-malicious byte-sequences or anomalous HTTP-requests. While this offers a direct view at the content of communications, it raises several problems: The computational overhead of DPI is large and scales directly with the amount of transferred data. Furthermore, DPI infringes user privacy when scanning the content of messages, potentially allowing the operator to access the vast amount of personal information sent over the internet. Finally, DPI is incapable of processing encrypted traffic, an issue that malicious actors increasingly exploit. Some firewall-providers have started to decrypt any transferred traffic in a man-in-the-middle manner, but the ethical legitimacy of such approaches is at the very least controversial.

Shallow packet inspection (SPI) in contrast only inspects packet headers and is therefore computationally less expensive as well as robust against encryption and to some degree application evolution. SPI is used both by rule- and signature-based systems such as Snort [8], as well as by an ever increasing number of machine-learning based methods that attempt to make broad generalisations about traffic solely from

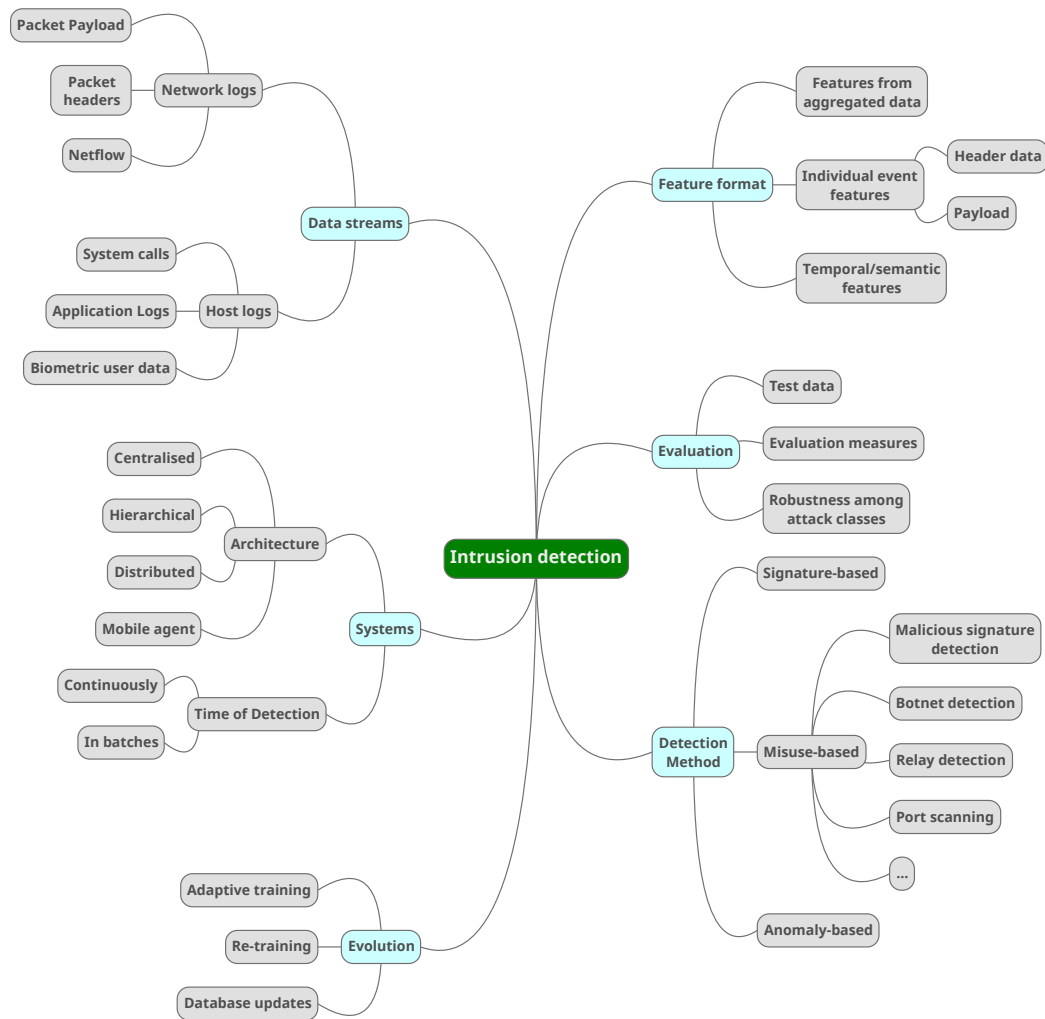


Figure 2.4: A broad overview over different aspects in an IDS

evaluating packet headers. These methods are discussed in more detail in the following chapters.

2.3.2 Misuse vs anomaly detection

Detection methods are the core of an IDS, and are therefore the most important design choice. Traditionally, two broad types of detection approaches are identified: a) misuse detection and b) anomaly detection.

Misuse detection aims at detecting a particular and well known reoccurring characteristic or pattern of a malicious behaviour. Two simple examples of such a characteristic are the large number of SYN packets sent by a host in a DoS attack, and the synchronised connection of many hosts to one server in a botnet. In misuse detection,

abnormal or malicious behaviour is therefore defined first before developing a model to distinguish the defined behaviour from other traffic.

In contrast, anomaly detection aims at building a model of normal system behaviour that is accurate enough to spot any malicious behaviour as traffic that deviates from the estimated model. Anomaly detection is principally more difficult than misuse detection since the traffic model has to incorporate potentially very heterogeneous traffic behaviours. By design, anomaly detection is more suitable to detect new and previously unseen malicious behaviour as it makes no definite assumptions on the anomalous behaviour. Misuse detection is robust against evolution of malware as long as defined malicious behaviours do not change.

In reality, anomaly and misuse detection are not necessarily mutually exclusive, and there is a fluent passage between the two. This is because many anomaly detection approaches choose a particular set of features to be modelled with a particular threat in mind. For instance, models for the number of connections of a machine are naturally suitable for detecting DoS attacks, port scans, or worm attacks.

2.3.3 Evolution of network anomaly-detection

In 1986, Denning and Neumann published the first IDS prototype model, called *Intrusion Detection Expert System* (IDES) [24]. The IDES model relies on point-based anomaly-detection based on both statistical and expert rules, and considers both network and host data. In the following years, several other systems were proposed that combine statistical and expert systems on security audit data, such as *MIDAS* [25] or *Haystack* [26]. In 1990, the first purely network-based anomaly detection methods emerged with the *Network Security Monitor*, and the *Network Anomaly Detection and Intrusion Reporter*. Both systems tried to profile user activity over time, and detect abnormalities with statistical estimates based on historic data. These systems were still relying on the notion that the aggregated values of service usage or destination hosts of an intruder is different enough from a legitimate user to be detected, and the corresponding detection resolution was only up to 5-minute intervals.

Between 2003 and 2005, several techniques to identify network-wide anomalies that are caused by worm propagation or scanning attacks were proposed first by Wagner and Plattner [27, 28], and followed by Anukool Lakhina [29, 30]. Lakhina's methods were more sophisticated, relying on principle component analysis and wavelet-modelling to detect sudden changes in the network, and were also the first anomaly-

based NID-methods with acceptable false-positive rates on real-world data.

In 1998, Martin Roesch developed *Snort*, the now most widely used NID-system [8]. Snort was the first system to identify individual packets as malicious, based on pre-defined rules on both packet content and packet header, but only detects misuse instead of anomalies. In 2002, both Kruegel et al. [31] and Mahoney and Chan [32] proposed packet-level anomaly-detection systems. Mahoney and Chan apply clustering techniques to both the packet header and the packet payload, while Kruegel et al. model the packet payload with Markov chains according to the specific service and request.

While Kruegel et al. only evaluate their model on proprietary DNS-data, Mahoney and Chan used the new DARPA 98 dataset, the first public benchmark NID dataset with labelled attack data [33]. This dataset, along with its derivatives, the KDD 99 and the NSL-KDD datasets [34, 35], were since then used heavily. In particular the format of the KDD 99 and NSL-KDD datasets, which consist of connection summaries with 41 features, influenced and simplified the way that NID-methods were designed. The overwhelming majority of systems proposed to detect access attacks since then have applied supervised or unsupervised ML-methods directly to these features, with each flow being classified independently from other traffic. Among the applied ML-techniques used for anomaly-detection are density-based models, support-vector machines (SVM), projection methods, Bayesian and fuzzy-logic classifiers, and neural networks such as autoencoders or deep belief networks. Prominent examples of this methodology include a hybrid SVM and neural clustering model trained with the DARPA 98 data by Shon and Moon in 2007 [36], and or an autoencoder model trained on the NSL-KDD dataset by Shone et al. in 2018 [37].

The existence of public benchmark datasets have enabled researchers to compare the performance of different methods, but also encouraged them to neglect operational considerations or structures akin to specific attacks and services. However, according to a recent study by Nisioti et al., several widespread evaluation malpractices which are discussed further in Section 5.3 decrease trustworthiness of published low volume access attack detection rates [13]. Overall, the academic progress of anomaly-based access attack detection has not progressed as much in the last decade, with false-positive rates still being magnitudes too high for operational deployment.

Further discussion of specific techniques related to anomaly detection using sequential models and to stepping-stone detection can be found in Sections 5.9 and 6.2.

2.4 Traffic datasets

In order to evaluate their ability to model the behaviour of hosts in a network and identify network intrusions, new methodologies have to be tested using existing datasets of network traffic. The dataset should ideally contain realistic and representative benign network traffic as well as a variety of different network intrusions. However, as network traffic contains a vast amount of information about a network and its users, it is notoriously difficult to release a comprehensive dataset without infringing the privacy rights of the network users. Furthermore, the identification of malicious traffic in network traces is not straightforward and often requires a significant amount of manual labelling work.

For that reason, only three datasets for network intrusion detection containing real world traffic exist: LANL-15/17, UGR-16, and LITNET-2020, all of which have lack however lack the necessary diversity of attack traffic. Significant efforts have been made to artificially create such datasets and thus bypass any privacy concerns. However, up to today, no artificial dataset truly resembles real network traffic in every aspect [13]. Table 2.1 at the end of this chapter provides an overview of the dimensions and characteristics of the most important NID datasets.

2.4.1 Real-world datasets

Los Alamos National Laboratory - Comprehensive, Multi-Source Cyber-Security Events 2015, and Unified Host and Network Data Set 2017 [38][39]

In 2015, the Los Alamos National Laboratory (LANL) released a large dataset containing **network flows** from their corporate computer network, which contains about 17,600 computers. The data was gathered over a period of 58 days with about 600 million events per day. The data only contains internal network connections, i.e. no flows going to or coming from computers outside the network are included. IPs and ports were de-identified, but are consistent throughout the data. Since the data stems exclusively from within the corporate network, it can be assumed that it shows more homogeneity in the observed traffic patterns than general network traffic.

The dataset furthermore contains a labelled set of red-team events which should resemble *pass-the-hash*-attacks. However, these events are not part of the network flow data and only contain information about the time of the attack and the attacked computer.

LANL released another dataset containing network flow traffic from their network in 2017 [40]. This dataset is similar to the one from 2015, but spans over a longer period of time, 90 days. Furthermore, it contains no labelled malicious activity.

UGR 2016 [41]

The UGR 2016 dataset was released by the University of Grenada and contains **network flow** data from a cloud service provider to a number of companies, and thus the data comes from a much less structured network than the LANL data. It contains both client's access to the internet and traffic from servers hosting a number of services. IP-adresses are consistently anonymised while network ports are unchanged. It is not ensured that all traffic coming from and going to a particular machine is captured. The dataset correspondingly covers a very long period, spanning from March to August of 2016, and containing about 14 GB of traffic per week.

The dataset contains labelled traffic from three real attacks, corresponding to IP-scanning and a spam mail campaign. The dataset also contains twelve days of synthetically generated data from a virtualised network along with controlled DoS and port scanning attacks and injected C&C traffic.

2.4.1.1 LITNET-2020 [42]

The LITNET-2020 from the Kaunas University of Technology Lithuania from 2020 was collected from an academic network over a timespan of ten months and contains annotated real-life benign and attack traffic. The corresponding network provides a large network topology with more than a million IP-addresses, and the data was collected in the form of network flows with more than 80 features. However, the dataset only contains traffic from high volume attacks such as DoS-, scanning, or worm attacks.

CAIDA and MAWIlab datasets [43, 44]

The *Center for Applied Internet Data Analysis* (CAIDA) started collecting network traces from a US high-speed backbone link in 2008, with the collection still ongoing. The data is available in anonymised yearly released datasets containing one hour of **packet headers**. Since the data is captured at a backbone, the capture is not necessarily complete, i.e. some packets in a connection could be routed via another backbone and not appear in the capture. It is furthermore not necessarily free from attack traffic.

Although this dataset has been used for intrusion detection before, it is more suitable for general internet traffic analysis.

Similarly to the CAIDA datasets, the MAWIlab dataset from the MAWI (Measurement and Analysis of the WIDE internet) research group contains **packet headers** from the WIDE backbone, with weekly updates. It is therefore similarly unstructured, anonymised, and not free from attack traffic. The dataset is labelled with the labels “anomalous”, “suspicious”, “notice”, and “benign” according to several heuristics, and has identified traffic from several worm and DoS attacks.

UNIBS 2009[45]

This now outdated dataset was collected on the campus network of the *University of Brescia* on three consecutive days in 2009. The dataset contains in total 79,000 anonymised TCP and UDP **network flows**, which is far smaller than the other datasets discussed here.

This dataset is not directed towards intrusion detection research, but was made as ground truth data for traffic classification. It therefore contains labels which indicate which of in total six applications generated the corresponding traffic flow.

2.4.2 Synthetic datasets

CICIDS 2017/2018 [46][47]

The CICIDS 2017 dataset, released by the *Canadian Institute for Cybersecurity* (CIC), contains five days of network traffic from 15 virtual machines. The network contains switches, routers, a web server, a modem, and a firewall in order to ensure a realistic network topology. The traffic data itself consists of **labelled benign and attack traffic**, and is available as 11 GB per day of **raw packets** with payloads, or as **network flows**.

The background traffic is not directly generated through user interactions on the machine, but by using a method to profile abstract user behaviour in different traffic protocol. However, it is not completely clear how much of the underlying structure of real traffic is lost in the process, and therefore how suitable this data is to build models of benign user activity.

The attack data of this dataset is one of the most diverse among NID datasets, as it contains attacks such as DoS, SQL-injections, Heartbleed, brute-forcing, probing, or XXS. These are not always successful in order to reflect actual attack scenarios.

However, the authors did not describe well how exactly the data from these attacks is generated.

CICIDS 2018: This dataset is generated in a similar fashion to the CICIDS-2017 dataset. The main differences are that the CICIDS 2018 data spans over three weeks and includes in total 450 hosts, but lacks the amount of web-attacks that is present in the CICIDS 2017 dataset.

UNSW-NB 2015 [48]

The dataset released by the *University of New South Wales* in 2015 contains traffic collected at the *Cyber Range Lab of the Australian Centre for Cyber Security*. The background traffic is generated from scripted activities on around 45 virtual machines, and is overlayed with replayed attack traffic using the *IXIA PerfectStorm tool*. The time span of the collection is in total 31 hours, and the data is available as **raw packets** and **network flows** along with two other data formats containing newly engineered features.

The attacks include a variety of DoS, reconnaissance, and access attacks. However, due to the synthetic injection of traffic from these attacks, it is unclear how close they are to real-world attack scenarios.

DARPA 1998 [49] & KDD Cup 1999 [34]

The *Defense Advanced Research Projects Agency* released the first major dataset to test network intrusion detection systems in 1998. The data stems from two experiments at the *MIT Lincoln Laboratory* where multiple victim hosts running Unix and Windows NT were subject of over 200 attacks of 58 different types. The data spans three weeks of training and two weeks of testing data and contains **raw packets** that are labelled.

The KDD Cup 1999 dataset was created by the *MIT Lincoln Laboratory* by processing portions of the 1998 DARPA dataset with new labels, and has been the most widely used dataset in network intrusion detection. It contains two million connections summaries in a new format similar to flow summaries, and in total 38 attack types.

The DARPA 98 dataset received a lot of criticism for its lack of realistic background traffic and the presence of artefacts from these simulations in the data. Naturally, as the KDD'99 data stems directly from the DARPA dataset, it also faces the same problems and criticism.

CTU 2013 [50, 51]

The *Stratosphere Laboratory* in Prague released this dataset in 2013 to study botnet detection. It consists of more than 10 million labelled **network flows** captured on lab machines for 13 different botnet attack scenarios. Additionally, the **raw packets** for the botnet activity is also available for attack analysis.

A criticism of this dataset is the unrealistically high amount of malicious traffic contained in the dataset, which makes it easier to spot it while reducing false positives. Furthermore, the way normal or background traffic is generated is described only poorly, and the labelling is imprecise.

2.4.3 Generative traffic models and traffic generation

Network traffic models and corresponding traffic generation models have been one of the cornerstones of network design and stress testing. The fidelity of specific traffic characteristics to real-world traffic depends on the purpose of the generator. Simple **performance measurement tools** such as *iPerf* or *Mausezahl* are used to explore network transmission reliability and efficiency, and therefore only send dummy packets that do not correspond to actual communication [52, 53]. The realism of network activity levels and spikes is a crucial aspect for these tools. The traffic is typically sent from a client device to a server device.

Traffic generators such as *IXIA PerfectStorm* aim to test the performance and scalability of **content-aware** firewalls and similar security devices [54]. They replay application traffic such as web traffic or video streaming as well as traffic from various attacks to test the reliability and overhead of a security device. The generated traffic therefore contains actual communication with semantically correct packet sequences and payloads, but offers little room to modify the transmitted traffic. The network set-up here is similar to tools like *iPerf*.

Traffic generation for NID-datasets is often performed using an arrangement of virtual machines that communicate to generate traffic, such the one used in the UGR-16 dataset as depicted in Fig. 2.5 [41]. The machines are configured to cover a range of network applications and protocols and communicate via scripted interactions. Attacks are typically performed from specifically designated attack machines. In Chapters 3 and 4, we discuss several shortcomings of this method.

Recently, traffic emulation tools that are powered through **Generative Adversarial Networks** (GANs) such as DoppelGANger [55] have attracted attention. These are

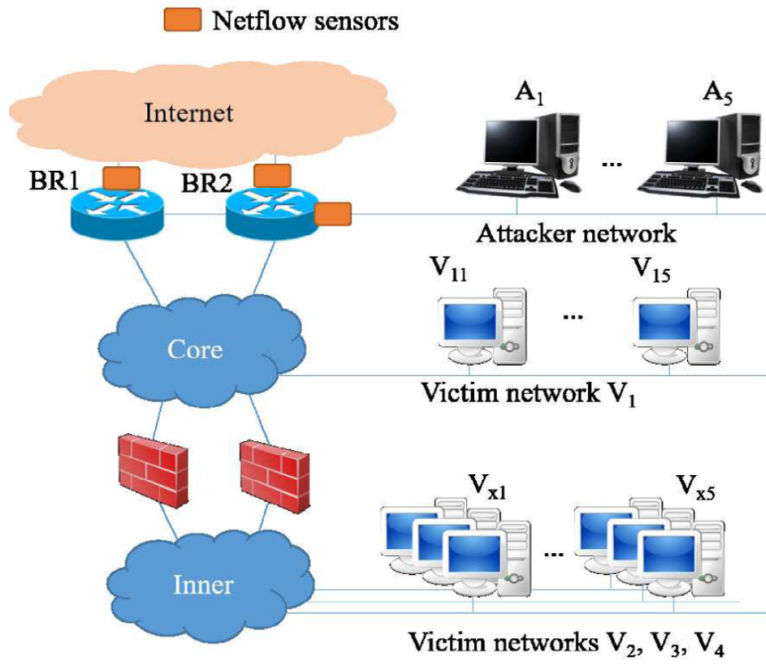


Figure 2.5: Setup of VM-machines in the UGR-16 dataset [41].

trained on real-world data and are capable of generating realistic sequences of traffic features. In contrast to actual traffic generators, no actual traffic is generated and the output of GAN-models so far has been limited to numerical features such as activity levels or port entropy distributions obtained after traffic processing. Approaches to generate sequences of packets or flows have so far not been successful [56, 57].

Dataset	Format	#hosts	Duration	Source	Attacks	Generation setup
LANL-15	Flows (only internal)	1760	58/90 days	enterprise network	Pass-the-hash	IP&Port-anonymisation
UGR-16	Flows	30	6 months	cloud provider (+12 days synthetic)	DoS, probing, C&C	IP-anonymisation, scripted activities&attacks
CAIDA	Packet headers	> 1 million	1h each	internet backbone	-	IP-anonymisation
MAWI	Packet headers	> 1 million	15 min each	internet backbone	DoS, worms, probing	IP-anonymisation
CICIDS-18	Pcap&ML-flows	450	3 weeks	VM-network	13 (DoS, probing, access attacks)	activity profiling, scripted activities&attacks
CICIDS-17	Pcap&ML-flows	30	5 days	VM-network	12 (DoS, probing, access attacks)	activity profiling, scripted activities&attacks
UNSW-NB-15	Pcap&ML-flows	45	31 hours	VM-network	9 (DoS, probing, fuzzing, decryption, access attacks)	scripted activities, replayed attacks

Table 2.1: Summary of the most commonly used NID-datasets.

Chapter 3

Requirements for Machine Learning

3.1 Introduction

Security oriented datasets describing computer networks are notoriously hard to obtain, and researchers struggle to evaluate new NID systems on suitable network traffic data. The lack of quantity, variability, meaningful labels, and ground truth has so far slowed scientific progress and objective and appropriate measurements on ML-based network security methods.

Privacy and security concerns discourage network administrators to release rich and realistic datasets for the public. Network traffic produced by individuals contains a mass of sensitive, personal information, such as passwords, email addresses, or usage habits, requiring researchers to expend effort anonymising the dataset [58]. To examine malicious behaviour, researchers are often forced to build artificial datasets using isolated virtual machines in a laboratory setting to avoid damaging operational devices. Background traffic is usually generated in real-time from scripts executed on the virtual machine, which constrains both the amount and heterogeneity of the data.

Existing network intrusion detection datasets are predominantly designed to support a broad range of applications, and are collected in a static manner, unable to be modified or expanded. This proves to be a serious defect as the ecosystem of intrusions is continually evolving. Furthermore, it prohibits a more detailed analysis of specific areas of network traffic due to the available data only being a fraction of the original dataset. To combat this, new datasets must be periodically built from scratch.

Allowing researchers to create datasets dynamically to circumvent these issues would be highly beneficial. Container networks have recently found adoption to conduct traffic generation experiments, such as by Fujdiak et al. [59] who use container-

ised web servers to generate DoS-traffic. The reason of using containerisation in these experiments is to provide a lightweight, scalable, and extendible framework that produces attack traffic in a secure way. In this chapter, we propose a framework, called **DetGen** that uses *Docker*-containers, that is designed to increase the control, monitoring, and experimental determinism in the traffic generation process to provide researchers with more information for model analysis [60]. Each Docker container is highly specialised in its purpose, generating traffic related to only a single application process. This shields the generation process from background activities and corresponding traffic events or disturbances during the collection. In combination with the scripting of a variety of well-defined and monitored communication *scenarios*, both for benign and attack traffic, and the controlled simulation of external effects, we can build a dataset with quasi perfect ground truth. Since many containerised applications are shared on the Docker Hub platform, implementation process is relatively easy.

The following results are discussed in this chapter:

1. We present a novel network traffic generation framework that is designed to improve several shortcomings of current datasets for NIDS evaluation. This framework is openly accessible for researchers on *GitHub* and allows for straightforward customisation.
2. We define four new requirements a network intrusion dataset should fulfil in order to be suitable to train machine-learning based intrusion detection methods.
3. We perform a number of experiments to demonstrate the suitability and utility of our framework.

This chapter is mostly consisting of work published in “Traffic generation using containerisation for machine learning” (Henry Clausen, Robert Flood, and David Aspinall, 2019 [1]).

3.1.1 Outline

The remainder of this chapter is organised as follows. Section 3.2 background information about application virtualisation methods. Section 3.3 discusses the problems that arise during the usage of existing NIDS datasets and concludes with a set of requirements we propose to improve the training and evaluation of machine-learning-based

methods. Section 3.4 describes the general design of our framework, and how it improves on the discussed problems in existing datasets. We also discuss a specific example in detail. Section 3.5 discusses several experiments to validate the improvements and utility our framework provides. Section 3.6 concludes the results and discusses limitations of our work and directions for future work.

3.1.2 Scope of DetGen and potential use-cases

The scope of DetGen is to generate traffic with near-deterministic control over factors that influence microscopic packet- and flow-level structures. DetGen separates program executions and traffic capture into distinct containerised environments to exclude any background traffic events, simulates influence factors such as network congestion, communication failures, data transfer size, content caching, or application implementation.

The purpose of DetGen is to better train and understand traffic models that model individual or short sequences of connections. Due to the controllable randomisation of influence factors and its extendibility, researchers can easily generate tailored traffic samples at scale to train and fine-tune their models on particular traffic types (an example of this is described in Section 3.5.3), or test the behaviour of models when encountering different data balances during training or evaluation to understand which structures the models fails on (a detailed description of this process is described in Chapter 4).

It is important to mention that since DetGen is currently only capable of generating atomic traffic events consisting of individual or short sequences of connections, the data generated with DetGen is not suitable to train or evaluate traffic models of network-wide features such as usage distributions or activity spikes as well as models that correlate events in a traffic stream over more than a few seconds. It is also not suitable to examine multi-step attacks yet as the included attack scenarios also consist of atomic interactions. The simulation of causal dependencies in attack sequences as well as for individual hosts is considered for future extensions of DetGen, as described in Section 3.6.2.

```
FROM ubuntu
MAINTAINER XYZ (email@domain.com)
RUN apt-get update
RUN apt-get install -y nginx
ENTRYPOINT ["/usr/sbin/nginx","-g","daemon off;"]
EXPOSE 80
```

Figure 3.1: Example of `Dockerfile` creating a nginx-container.

3.2 Background

3.2.1 Containerisation with Docker

Virtual machines (VMs) share the same hardware infrastructure as the host machine. VMs necessitate the use of hypervisors, software responsible for sharing the host OS's hardware resources, such as memory, storage and networking capabilities. OS-level virtualisation, also known as *containerisation*, is a virtualisation paradigm that has become popular in recent years due to its lightweight nature and speed of deployment. In contrast with standard VMs, containers forego a hypervisor and the shared resources are instead kernel artefacts, which can be shared simultaneously across several containers. Although this prevents the host environment from running different operating systems, containerisation incurs minimal CPU, memory, and networking overhead whilst maintaining a great deal of isolation [61].

The main advantage of using containers for traffic generation is the isolation of individual applications. This enables us to gather ground truth about the traffic origin, and enables us to easily extend, modify, and scale our traffic generation framework, which would not be possible when relying on VMs.

Docker container *Docker* is a software platform that allows for the creation, maintenance and deployment of containers. In Docker's terminology, a container is a single, running instance of a Docker *image*. Docker images are defined via a text file known as the `Dockerfile`, which consists a series of commands that modify an underlying *base image*, usually a containerised OS. Example commands include installing libraries and copying files. Figure 3.1 displays a simple example of `Dockerfile`.

After each command is executed, the intermediate, read-only image is saved as a *layer*. These layers can be shared between containers. When a Docker image is run as a container, a final read-write layer is added and when the container is later stopped,

this layer is discarded, preserving the integrity of the underlying layers. This allows Docker containers to be run repeatedly whilst always starting from an identical state.

Individual Docker containers are intended to be highly specialised in their purpose with each container running only a specific piece of software or application. Commonly used base images — such as Alpine Linux — have minimal background processes running during a container’s lifetime. This means that the network traces of a Docker container can be associated with a specific application. The one-to-one correlation between containers and network traces allows us to produce labelled datasets with fully granular ground truths.

The Docker software platform includes a cloud-based repository called the Docker Hub [62] which allows users to download and build open source images on their local computers. At the time of writing, nearly 2.5 million images are available from the Docker Hub. Some common software — such as popular webserver and databases — have officially maintained images. We use these as far as possible to simplify the production of our scenarios, and keep them close to software configuration used in practice.

Docker Networking Docker allows the creation of virtualised networks with one or more subnetworks, to which containers can connect via a virtualised network bridge. Containers attached to the bridge network are assigned an IP address and are able to communicate with other containers on their subnetwork. Containers can furthermore be connected to a host network, which allows communication with external networks using NAT via the host interface.

To host containers in an isolated network, we can create our own user-defined bridge networks, which provides greater isolation between containers [63]. Furthermore, this allows us to fix the subnet and gateway for our networks as well as the IP addresses of our containers, which simplifies scripting scenarios. Docker allows containers to share the same network interface. This enables us to assign the same IP address to multiple containers.

NetEm The Docker engine also provides network access to Linux traffic control facilities such as *NetEm*, on which we will rely in this project. NetEm is a Linux toolkit for testing protocols by emulating properties of wide area networks [64]. It allows the user to emulate variable delay, loss, duplication and re-ordering of packets on particular network interfaces.

```
version: '3'
services:
  webserver:
    image: nginx:alpine
    ports:
      - "80:80"
    networks:
      - app-network
  db:
    image: mysql:5.7.22
    ports:
      - "3306:3306"
    environment:
      MYSQLDATABASE: laravel
      MYSQLROOTPASSWORD: your_mysql_root_password
    networks:
      - app-network
networks:
  app-network:
    driver: bridge
```

Figure 3.2: Example of Docker-compose file launching a nginx- and a mysql-container in an isolated network.

Docker Compose Applications built using the Docker framework often need more than one container to operate, for example an Apache server and a MySQL server running in separate containers. We must build and deploy several interconnected containers simultaneously. Docker provides this functionality via `Docker compose`, a tool that allows users to define the services of multiple containers as well as the properties of virtual networks in a YAML file. By default, this file is named *docker-compose.yml*. This allows for numerous containers to be started, stopped and rebuilt with a single command in a consistent manner. This is particularly significant for our purposes; with `Docker compose`, we can launch several containers in a specific order, with a specific network configuration, whilst running specific commands within each container on start up. This ensures that our interactions are deterministic, barring any added randomisation. Figure 3.2 displays an example of a simple `Docker compose` file.

3.3 Problems in modern datasets

The difficulty of obtaining malicious traffic in real-world captures means that the performance of new network intrusion detection algorithms are almost exclusively evaluated on synthetic datasets. Potential disadvantages of synthetic compared to real-world datasets have been discussed by several authors [14, 65]. However, none address problems in the particular design of such synthetic testbeds that are holding machine-learning based methods back in performance and from getting more widespread application. Here, we focused on this aspect and four design problems common among modern synthetic datasets.

Lack of variation To generate benign traffic, a selection of activities is scripted and executed on virtual machines. Activities are selected to cover the most prominent protocols, but seldom to cover the range of subactivities that each protocol offers. Instead, the manner in which each protocol is used is highly-restricted, and there are doubts about whether this traffic is representative of its real-world equivalent usage [14]. An illustrative example of the restricted protocol activity in synthetic datasets can be seen in the CIC-IDS 2017 dataset. Here, the vast majority of successful FTP transfers consist of a client downloading a single text file containing the Wikipedia page for ‘Encryption’ several hundred times in a day. In reality, FTP is used for a large number of tasks, which can occur in random order with varying input sizes and parameters.

In addition to that, implemented test bed environments are usually separated from external influence or even virtualised, which isolates them from fluctuations and faults introduced by the complexity of modern networking. These include packet delays through network congestions, unexpected connection drops or resets, and out-of-order arrivals, all of which lead to variations in the response behaviour of particular services.

This general lack of variation in individual protocols leads to observed homogeneity both on a packet exchange level and on a network flow level, and thus to clearer structures in the data. Identifying separations of malicious and benign activity or between different services consequently becomes easier, which leads to overoptimistic results in the evaluation of machine-learning based methods. It is therefore clear that traffic variation is a crucial aspect of a comprehensive intrusion detection dataset. parameter

Lack of ground truth To evaluate machine-learning-based methods that distinguish between different types of network traffic data, we need to verify that separating structures in the model correspond to distinct computational actions, using traffic labels. Most obvious is the labelling of benign and malicious traffic. More granular labels are desirable to distinguish between several different types of network traffic. An example for this is the design of ‘stepping stone’ detection methods, where researchers try to detect connections relayed over a jump-host. Similarity or correlation metrics that measure the closeness of two connections are a popular tool. To understand such a measure, ground truth about how computationally similar two connections are and what type of behaviour they represent is necessary. Other areas that look at small-scale traffic structures and would benefit from detailed traffic labels include protocol verification, traffic classification, traffic disaggregation, or exploit discovery.

Ground truth labels for network traffic are hard to obtain. The network traffic produced by a typical PC will invariably contain traffic originating from background processes, such as software updates, authentication traffic, network discovery services, advertising features, as well as many other sources. To separate traffic from different origins retrospectively is often hard, if not impossible. Source attribution through port numbers is unreliable because port numbers can be dynamically allocated and are not restricted to particular processes, and processes can open connections on multiple ports at the same time. All of these reasons mean that the identification of different computational operations from captured traffic is often infeasible. Therefore, no public NID dataset currently considers the inclusion of ground truth traffic labels.

Static design A released dataset can only contain data that is representative a system at the time of creation. In contrast to other many other data sources for machine-learning, network traffic, both benign and malicious, is constantly changing as computational protocols and systems evolve. All available NIDS datasets today have been created in a static manner, so that a fixed test bed of host machines is created designed to contain specific vulnerabilities to the selected attacks. This makes it very hard to change the test bed and thus adjust the dataset to updated traffic structures. Allix et al. [66] claim that it is impossible to release a NID dataset that is truly representative of the real-world attacks due to the inherent secrecy of the intrusion ecosystem and the rate at which it develops.

Limited size Today's machine-learning revolution was supercharged by the exponential growth of available data. Larger amounts of data mean that a given model can identify more complex structures that remain invariant in noisy environments and thus generalise better. Although the amount of globally transmitted network traffic is growing every year, the size of available NIDS datasets is limited by small host numbers, typically 5-10, and short capture periods, at maximum a 5-6 weeks, inherent to test bed captures. This means that traffic models can experience difficulties to generalise over specific traffic types which represent a smaller fraction of the total dataset. In an ideal setting, researchers would have the ability to generate arbitrary amounts of specific traffic types.

3.3.1 Dataset Requirements

The primary task of this project is to provide a suite of Docker container compositions that is capable of generating traffic datasets suitable for machine-learning-based intrusion detection systems. This container suite is designed to address the criticism of current NIDS datasets discussed in Section 3.3. For this, we created a set of requirements that a modern intrusion detection dataset has to fulfil to address the problems discussed in Section 3.3:

Variation To ensure that we produce representative data for modelling, we want the traffic generated by our container suite to cover a sufficient number of protocols that are commonly found in real-world traffic and existing datasets. For malicious traffic, we want to ensure that the attacks are modern and varied, both in purpose and in network footprint. For each protocol, we want to establish several capture scenarios to encompass the breadth of that protocol's possible network traces. Communication between containers should be subject to the same disturbances and delays as in a real-world setting.

Ground truth Since ground truth is a main focus of this work, we want a capture scenarios to be consistent and reproducible in the traffic they generate. This way, we can be certain that a particular traffic trace corresponds to the capture scenario it was generated by, and can thus relate individual traffic events to computational operations. We discuss what it means for a scenario to be reproducible in detail in Section 3.5.1.

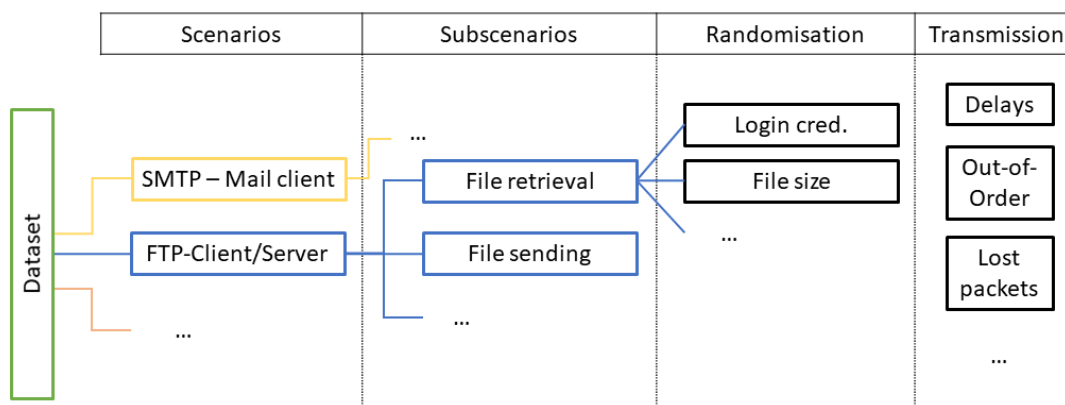


Figure 3.3: Visualisation of the different levels at which traffic variation is introduced in DetGen.

Modularity Traffic capture scenarios should be implemented in a modular way allow for a straightforward addition or modification of traffic capture modules without disrupting the rest of the container suite. This reduces the effort to adjust a dataset to changing traffic patterns and allows the addition of modern attacks traffic.

Scalability Each capture scenario should be running in a scalable manner to allow generation of large data quantities.

3.4 Design

To cover a range of activities, the containers in our framework are arranged in different configurations corresponding to particular *capture scenarios*. Running a given capture scenario triggers the launch of several Docker containers, each with a scripted task specific to that capture scenario. A simple exemplary capture scenario may consist of a containerised client pinging a containerised server. We ensure that each Docker container involved in producing or receiving traffic will be partnered with a `tcpdump` container, allowing us to collect the resulting network traffic from each container's perspective automatically.

We outline different stages within the creation of a dataset at which traffic variation is introduced. Figure 3.3 visualises this process.

3.4.1 Scenarios

We define a *scenario* as a series of Docker containers interacting with one another whereby all resulting network traffic is captured from each container's perspective. This constructs network datasets with total interaction capture, as described by Shiravi et al. [67]. Each scenario produces traffic from either a protocol, application or a series thereof. Both benign and malicious activities are implemented as scenarios. Examples may include an FTP interaction, a music streaming application and client, an online login form paired with an SQL database, or a C&C server communicating with an open backdoor. A full list of currently implemented scenarios can be found in Section 3.4.9.

Each scenario is designed to be easily started via a single script that allows the user to set the length of the capture time, and the specification of particular subscenarios, discussed below. Scenarios can be repeated indefinitely without further instructions and be run in parallel, therefore allowing the generation of large amounts of data.

Our framework is modular, so that individual scenarios are configured, stored, and launched independently. Adding or reconfiguring a scenario has no effect on the remaining framework.

3.4.2 Subscenarios

In contrast to scenarios, *subscenarios* provide a finer grain of control over the traffic to be generated, allowing the user to specify the manner in which a scenario should develop. The aim of having multiple subscenarios for each scenario is to explore the full breadth of a protocol or application's possible traffic behaviour. For instance, the SSH protocol can be used to access the servers console, to retrieve or send files, or for port forwarding, all of which may or may not be successful. It is therefore appropriate to script multiple subscenarios that cover this range of tasks.

The same applies to malicious activity. For instance, it would be naive for an SSH password bruteforcing scenario to always successfully guess a user's password. Instead, we include a second subscenario in which the password bruteforcer fails.

Subscenarios are specific to particular scenarios and can be specified when launching that scenario.

3.4.3 Randomisation within Subscenarios

Scripting activities that are otherwise conducted by human operators often leads to a loss of random variation that is normally inherent to the activity. As mentioned in Section 3.3, the majority of successful FTP transfers in the CIC-IDS 2017 data consist of a client downloading a single text file. In reality, file sizes, log-in credentials, and many other variables included in an activity are more or less drawn randomly, which naturally influences traffic quantities such as packet sizes or numbers.

To account for these fluctuations, we identify variable input parameters within scenarios and their subscenarios and systematically draw them randomly from a suitable distribution. Passwords and usernames, for instance, are generated as a random sequence of letters with a length drawn from a Cauchy distribution, before they are passed to the corresponding container. Files to be transmitted are selected at random from a larger set of files, covering different sizes and file names.

3.4.4 Network transmission

Docker communication takes place over virtual bridge networks, so the throughput is far higher and more reliable than in real-world networks, with the Docker virtual network achieving a bandwidth of over 90 Gbits/s when measured using iPerf [68]. This level of speed and consistency is worrying for our purposes as packet timings will be largely identical on repeated runs of a scenario and any collected data could be overly homogeneous.

To retard the quality of the Docker network to realistic levels, we rely on emulation tools. As discussed in section 3.2.1, NetEm is a Linux command line tool that allows users to artificially simulate network conditions such as high latency, low bandwidth or packet corruption in a flexible manner.

Although it is relatively straightforward to apply NetEm commands to a Docker Bridge network, we decided not to invoke NetEm in this manner as this would cause all network settings of all containers to be identical, such as all containers in a scenario having a latency of 50ms. Instead, we developed a wrapping script that applies NetEm commands to the network interface of a given container, providing us with the flexibility to set each container's network settings uniquely. This script randomises the values of each parameter, such as packet drop rate, bandwidth limit, latency, ensuring that every run of a scenario has some degree of network randomisation if desired.

3.4.5 Capture

To capture traffic, we use containers running `tcpdump`, a widespread and free packet analyser software that can capture packets arriving at or leaving from a network interface [69]. We attach `tcpdump` containers on every interface in the virtualised docker network and write packets into separate capture files. This allows us to capture traffic from the perspective of every container in a scenario, giving a complete view. In Section 3.4.8, we discuss how the collected capture files are coalesced into one dataset.

3.4.6 Implementation Process

The implementation process for each scenario follows broadly the same outline:

1. Select containers which provide the required services and identify the *primary* container/s for a given scenario which is/are dictating the container interaction. Then create and build a Dockerfile containing all necessary dependencies.
2. Identify different ways to use the service of the given scenario and define them into a set of subscenarios.
3. Design and implement the behaviour for secondary containers to provide the required service to the primary container(s).
4. For each subscenario, identify variable input values and their appropriate range; then systematically implement their generation from appropriate distributions covering this range.
5. Add `tcpdump` containers to every network interface.
6. Create a `Docker compose` file that launches all containers simultaneously.
7. Finally, write a script that, upon running, calls this `Docker compose` file, applies a network emulation script to each container network interface, and allows the user to specify how long and how many times a scenario should be run.

Following the Docker guidelines [70], each container in our framework consists of a single service with a specialised purpose, with as few additional dependencies as possible. Moreover, we ensure that there are minimal inter-dependencies between the containers of a scenario. This allows us to easily modify and update containers as new versions of the underlying software are released.

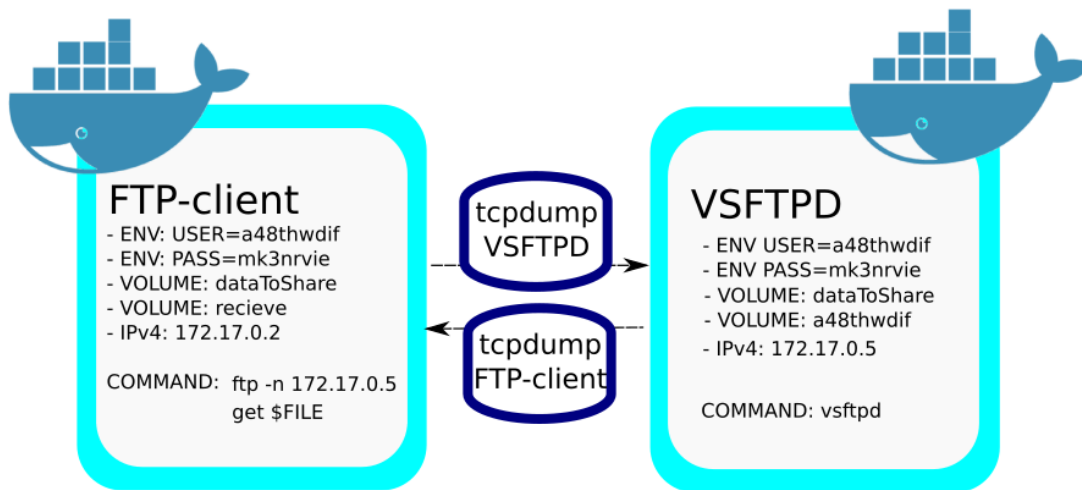


Figure 3.4: Diagram of FTP scenario

3.4.7 Simple Example Scenario - FTP server

We review the design of a prototypical capture scenario, namely, an FTP server and client interaction. The interaction is initiated by a single script, which allows the user to specify the length of the interaction, the number of times the interaction takes place as well as the specific subscenario. The script generates a random ftp username and password, creating the necessary *User* directory on the host machine before calling the Docker-compose file which creates a bridge network. Subsequently, the necessary containers are then started which, in this case, consist of a VSFTPD server, a client with ftp installed and two containers running tcpdump to capture all of the traffic emitted and received by the client and server respectively into separate .pcap-files. These .pcap-files are shared with the host machine via a shared volume. The host machine also shares:

- A *dataToShare* volume containing files that can be downloaded by the client.
- The *User* directory with the server, which contains the same files as the *dataToShare* folder.
- An empty *receive* folder with the client into which the files will be downloaded.
- The random username and password is shared with the client container so it can authenticate itself to the server.

Up to this point, no network traffic has been generated and the containers are now ready to begin communicating with one another. For this particular interaction between

an FTP server and client, we want to ensure that it is possible to capture the many ways in which an FTP session may develop. For instance, the client may seek to download files via the `get` command or the `put` command, alongside many other possibilities. We define 13 possible capture subscenarios intended to encapsulate a wide range of potential FTP sessions. These include downloading a single file using `get` or `put`, downloading every file using `mget` or `mput`, deleting every file from the server and requesting files from the server without the necessary authentication.

After the scenario ends, both the *User* directory and any downloaded files are removed from the host machine. The containers are then stopped and the bridge network is torn down. All necessary containers, volumes and scripts are in the same position prior to initiating the scenario — barring any generated `.pcap`-files — allowing for the scenarios to be started repeatedly with minimal human interaction. The `.pcap`-files are tagged with information about the time of creation, executed scenario and subscenario, and the container generating the traffic.

3.4.8 Dataset creation

Our framework generates network datasets consisting of a single interaction, but it is possible to coalesce these datasets to create larger datasets with a wide variety of traffic, albeit with some caveats. Due to the networking constraints of the Docker virtual network, such as limitations regarding clashing ports, running many of our Docker scenarios simultaneously over a large period of time is infeasible. Thus, to ensure that the generated traffic is suitably heterogeneous, numerous datasets must be generated before being coalesced into a main dataset. If done naively, this presents a problem. As discussed by Shiravi et al. [67], merging distinct network data in an overlapping manner can introduce inconsistencies. For instance, if one wanted to create a dataset containing both normal webserver traffic and traffic originating from a Denial of Service attack, it would not work to generate these two datasets separately before merging them together. If these two events really did occur simultaneously, the high network throughput of the latter would likely effect the packet timings of the former.

To avoid such inconsistencies, we create larger datasets by collecting data in consecutive chunks of fixed time. Within each chunk, several scenarios are run simultaneously. All `.pcap`-files collected during a given chunk can be merged together. It is then simple to stitch together all of these chunks into a single `.pcap`-file using a

combination of Mergecap [71] and Editcap [72]. This allows us to shift the timings of each .pcap-file by a fixed amount such that all of our chunks occur in succession whilst maintaining the internal consistency of each chunk.

3.4.9 Implemented scenarios

Our framework contains 29 scenarios, each simulating a different benign or malicious interaction. The protocols underlying benign scenarios were chosen based on their prevalence in existing network traffic datasets. These datasets consist of common internet protocols such as HTTP, SSL, DNS, and SSH. According to our evaluation, our scenarios can generate datasets containing the protocols that make up at least 87.8% (MAWI), 98.3% (CIC-IDS 2017), 65.6% (UNSW NB15), and 94.5% (ISCX Botnet) of network flows in the respective dataset. Our evaluation shows that some protocols that make up a substantial amount of real-world traffic are glaringly omitted by current synthetic datasets, such as BitTorrent or video streaming protocols, which we decided to include.

In total, we produced 17 benign scenarios, each related to a specific protocol or application. Further scenarios can be added in the future, and we do not claim that the current list exhaustive. Most of these benign scenarios also contain many subscenarios where applicable.

The remaining 12 scenarios generate traffic caused by malicious behaviour. These scenarios cover a wide variety of major attack classes including DoS, Botnet, Brute-forcing, Data Exfiltration, Web Attacks, Remote Code Execution, Stepping Stones, and Cryptojacking. Scenarios such as stepping stone behaviour or Cryptojacking previously had no available datasets for study despite need from academic and industrial researchers.

We provide a complete list of implemented scenarios in Table 3.1.

3.5 Validation experiments

A framework that generates network traffic does not necessarily provide realistic and useful data. To evaluate the utility of our Docker framework, we construct a series of experiments. We have two goals in mind. First, we want to demonstrate that the traffic generated is sufficiently representative of real-world traffic. Second, we want to demonstrate that having a framework to continually generate data compared to static

Name	Description	#Ssc.
Ping	Client pinging DNS server	1
Nginx	Client accessing Nginx server	2
Apache	Client accessing Apache server	2
SSH	Client communicating with SSHD server	5
VSFTPD	Client communicating with VSFTPD server	12
Scrapy	Client scraping website	1
Wordpress	Client accessing Wordpress site	1
Syncthing	Clients synchronise files via Syncthing	1
Mailx	Mailx instance sending emails over SMTP	2
IRC	Clients communicate via IRCd	2
BitTorrent	Download and seed torrents	3
SQL	Apache with MySQL	2
NTP	NTP client	2
Mopidy	Music Streaming	5
RTMP	Video Streaming Server	1
WAN Wget	Download websites	5
SSH B.force	Bruteforcing a password over SSH	3
URL Fuzz	Bruteforcing URL	1
Basic B.force	Bruteforcing Basic Authentication	2
Goldeneye	DoS attack on Web Server	1
Slowhttptest	DoS attack on Web Server	4
Mirai	Mirai botnet DDoS	3
Heartbleed	Heartbleed exploit	1
Ares	Backdoored Server	3
Cryptojacking	Cryptomining malware	1
XXE	External XML Entity	3
SQLi	SQL injection attack	2
Stepstone	Relayed traffic using SSH-tunnels	2

Table 3.1: Currently implemented traffic scenarios along with the number of implemented subscenarios

datasets benefits evaluating the efficacy of intrusion detection systems.

The first experiment provides a general verification of the reproducibility of our framework, which is required for guarantee the ground truth of the produced data. The second experiment demonstrates that the WAN-characteristics we emulate for our

data make it quasi non-distinguishable from real WAN traffic. Our third experiment then demonstrates the advantage of unlimited data generation capabilities for training ML-based traffic classification.

3.5.1 Reproducible scenarios

To provide ground truth, we have to guarantee that our implemented scenarios and subscenarios are consistent and reproducible upon repeated execution. This applies both to consistency for external influences on the host, such as increased computational load, as well as internal consistency of the implemented script execution.

It is impossible to guarantee that each scenario will produce a truly ‘deterministic’, or repeatable, output due to differences in network conditions, computational times, or input. Instead, we aim for our data to be *reproducible up to networking and computational differences*. This means that when running a scenario multiple times, we expect the quantities of most packets to be largely identical. We do expect some packets to exhibit greater variation due to non-determinism in the underlying protocols, Fig. 3.5 outlines this behaviour in terms of interarrival times and packet sizes.

To measure how consistent our scenarios are, we generate 500 .pcap files for three different implemented scenarios, namely the Apache, the VSFTPD, and the SSH scenario. These were generated consecutively under different host CPU load. We did not apply any delays or other NetEm traffic controls.

We assess the consistency of a scenario across different .pcap files by comparing all generated .pcap files pairwise. We measure this by the similarity of the connections captured.

To test the similarity of two connections, we extract the sample distributions of the packet interarrival times and packet sizes overall, upstream and downstream. We define two connections as similar if the two distributions for each of these quantities pass an equality test. We use the two-sample Kolmogorov-Smirnov (K-S) test, a non-parametric statistical test for the equality of two continuous one-dimensional distributions [73], with a p-value of 0.01.

As all tested files passed this similarity test, we conclude that these scenarios yield consistent and reproducible results. As other scenarios follow the same setup and launch commands, we expect the results to stay the same as long as the involved containers are consistent in their behaviour.

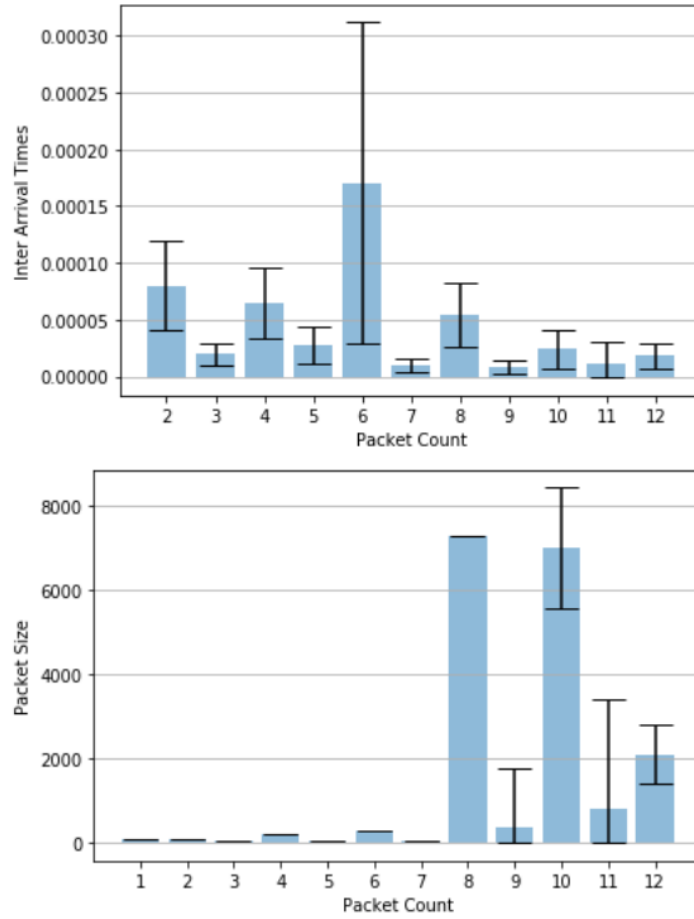


Figure 3.5: Means of IATs & packet sizes along with standard deviation bars for the first twelve packets in the Apache scenario.

3.5.2 Exploring Artificial Delays

Most traffic our framework generates is transported over Docker’s virtual network and therefore does not succumb to problems associated with normal network congestion, such as packet loss, corruption and packets arriving out of order. A realistic dataset should include these phenomena, which is why we developed wrapping scripts that allow us to artificially add delays as well as packet loss and corruption, using NetEm. Choosing the parameters is not straightforward; it is not clear how close to real-world traffic such network emulation techniques are. This is especially true for packet delays, which are described by continuous distributions and often have temporal correlation.

Furthermore, the high effective bandwidth of the Docker virtual network resulted in traffic with extremely short inter-arrival times (IATs, defined as the time between two packet arrivals). Therefore, we devote considerable time to demonstrating that it is possible for traffic generated by our Docker framework to conform to real-world IAT

distributions when altered using NetEm.

3.5.2.1 Datasets

We create two classes of datasets, one which is representative of ‘real-world’ traffic, and one which has been generated from our Docker framework. For simplicity, we only consider datasets consisting of FTP traffic.

For the real-world dataset, we set up a containerised VSFTPD server running on a *Google Compute* virtual machine located in the Eastern United States, and a containerised FTP client on our local host. We then ran a series of our scripted interactions between the two machines, generating 834 megabytes of data in 250964 packets. These interactions consisted of several FTP commands with various network footprints. We collect all data transmitted on both the server and the client. We call this data the *Non-Local* dataset.

We then repeat this process using the same container setup, but across the Docker virtual network on a local machine. We repeat this process several times, generating several *Local* datasets under a variety of emulated network conditions, discussed in Section 3.5.2.2. Our *Local* datasets vary slightly in size, but are all roughly 800 megabytes with 245000 packets.

3.5.2.2 Methodology

NetEm allows us to introduce packet delays according to a variety of distributions, namely uniform, normal, Pareto and Paretonormal¹. Furthermore, NetEm adds delays according to modifiable distribution tables, and so it is trivial for us to add a Weibull distribution, which along with Pareto distributions have been shown to closely model packet IATs [74, 75]. In total, we test the efficacy of four distributions to model inter-arrival times — normal, Pareto, Paretonormal and Weibull.

We generate several *Local* datasets by delaying traffic according these distributions, performing an exhaustive grid search over their means and standard deviations. Initial experiments revealed that introducing delays with a mean in the range of 40 ms to 70 ms produced the best results. Setting the jitter of the distribution too high resulted in the repeated arrival of packets out of order, therefore we further limit the grid search to

¹This Paretonormal distribution is defined by the random variable $Z = 0.25 * X + 0.75 * Y$, where X is a random variable drawn from a normal distribution and Y is a random variable drawn from a Pareto distribution.

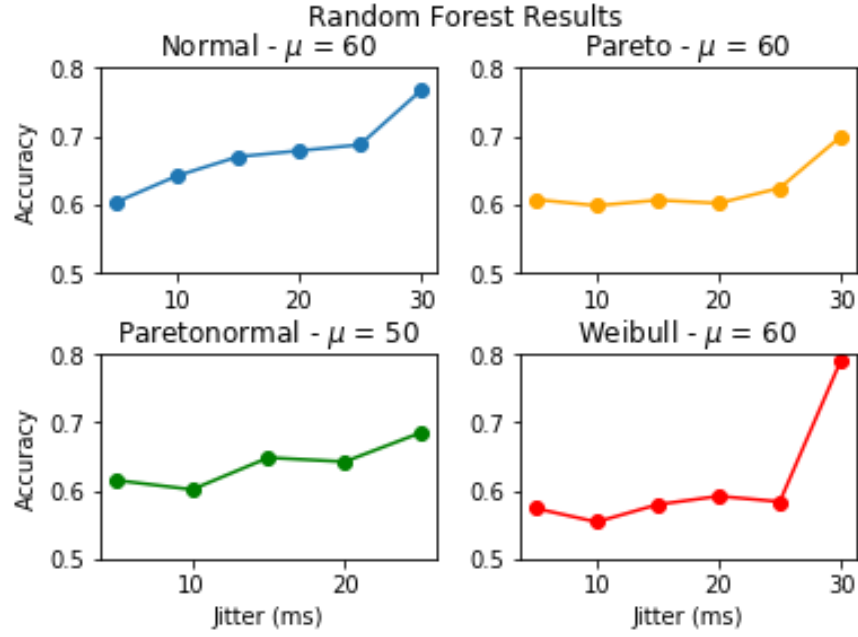


Figure 3.6: Results of Random Forest Classifier for a given distribution at the best performing delay mean μ . Note that a score of .5 indicates total indistinguishability.

jitter values in 5 ms intervals up to half of the value of the mean. In total, we generate 88 *Local* datasets.

Our goal is to discover the *Local* dataset whose packet timings most closely resemble those of our *Non-Local* dataset. To do this, we extract the IATs and packet sizes from our datasets on a packet-by-packet basis and store these results in arrays. We measure the similarity between two of these arrays by training a Random Forest classifier to distinguish between them. We say that if the Random Forest correctly classifies each packet with a success rate of only 50% then it is no better than randomly guessing and, as such, the inter-arrival times of these two arrays are indistinguishable from one another for the Random Forest.

To perform this measurement, we concatenate one *Local* dataset array with our *Non-Local* dataset array, label the entries and then shuffle the rows. We proportion this data into a training set and a testing set using an 80-20 split. We then feed this training data into a Random Forest with 1000 trees and fixed seed, and then record the accuracy of this Random Forest on the test set. We repeat this process for every single *Local* dataset.

3.5.2.3 Results

Table 3.2 summarises the values of the mean and jitter for a given distribution that produced the worst results from the random forest classifier.

DISTRIBUTION	MEAN	JITTER	RF ACCURACY
NO DELAYS (BASELINE)	0	0MS	0.8176
CONSTANT DELAY	40MS	0MS	0.6730
NORMAL	60MS	5MS	0.6028
PARETO	60MS	10MS	0.5979
PARETONORMAL	50MS	10MS	0.6015
WEIBULL	60MS	10MS	0.5540

Table 3.2: Worst Random Forest accuracy rates for a given distribution

To establish a baseline, we compare the traffic generated from our Docker scenario to that of the *Google Compute* data with no added delays. In this case, the Random Forest was able to distinguish between the two datasets, achieving an accuracy of over 90%. The classification accuracy is worsened considerably by introducing network delays, with the best results being achieved using a Weibull distribution with a mean of 60 ms and a jitter of 10 ms, leading to an accuracy of just 55%. Results for Pareto and Weibull distributions seem to yield consistent results for differing jitter values. Although not completely indistinguishable, this proves that using NetEm we can emulate WAN properties very closely.

3.5.3 Advantages of Dynamic Dataset Generation

Having examined whether our Docker framework is capable of emulating real-world IATs, we explore their utility in traffic classification to demonstrate the advantages that our framework provides compared to static, unlabelled datasets.

Machine-learning techniques are a popular tool for traffic classification, with many successful published classifiers. Furthermore, inter-packet arrival times have been shown to be a discriminative feature [76, 77]. However, these methods considered datasets consisting of completed traffic flows, limiting their use in, say, a stateful packet inspector. On-the-fly classifiers are also successful. Jaber et al. [78] showed that a K-means classifier can classify flows in real-time solely based on IATs with precision exceeding 90% for most protocols within 18 packets. Similarly, Bernaille et al. [79]

demonstrated that a K-means classifier can precisely classify traffic within five packets using only packet size as a feature.

However, Jaber et al. [78] only evaluated their traffic classifier with training and testing data drawn from the same dataset containing traces of a single network; there is no measure of how this model may generalise to other networks with differing conditions. Furthermore, they were limited to using unsupervised machine learning algorithms to classify their traffic as their datasets had no ground truth.

We attempt to replicate these results within our Docker framework with some adjustments. As we can generate a fully accurate ground truth, we attempt to segregate application flows based on their packet IATs using supervised learning techniques. Moreover, we then measure this model's ability to generalise by expanding our dataset to include traffic from networks with differing bandwidth and latency.

3.5.3.1 Data & Preprocessing

Our goal is to measure a classifier's ability to generalise across datasets. Therefore we construct two datasets using our Docker framework, both containing the same number of network traces from the same containers.

For our first dataset, we generate .pcap-files, each containing traffic from one of 16 different classes: HTTP (Client & Server), HTTPS (Client & Server), RTMP (Client, Server & Viewer), SSH (Client & Server), FTP (Client & Server), IRC (Client & Server), SMTP, SQLi and DoS traffic. To prevent class imbalance, we generate 200 .pcap-files for each of the 16 classes, resulting in 3200 total files. To more accurately emulate potential network conditions, we use our NetEm scripts to apply a unique delay to every container involved in a scenario. These delays follow a Pareto distribution with random mean between 0 and 100 milliseconds and random jitter between 0 and 10 milliseconds. We then preprocess this data by removing all but the first 12 packets of each .pcap-file. We extract the 11 inter-arrival times separating the 12 packets, which act as our feature vectors. We collect these feature vectors for each class along with a class label, and store collected feature vectors from all 3200 .pcap-files in a 12 x 3200 array. We call this our *Primary* dataset.

We then repeat this process to generate a second dataset, changing the properties of our emulated network. Again, we delay all traffic using a Pareto distribution, however, this time we select a random mean in the range of 100 to 500 milliseconds and random jitter between 0 and 50 milliseconds. The subsequent preprocessing of our data remains unchanged. We call this our *Secondary* dataset.

3.5.3.2 Methodology

First, we attempt to reproduce the results presented by Jaber et al. [78] by training a Random Forest with 100 trees to classify application flows based on packet IATs. We do this by proportioning our Primary dataset into training and testing sets using an 80-20 split. We then train and test our Random Forest repeatedly, first considering the classification accuracy based on the IATs of only the first two packets, then the first three packets and so on, up to 12 packets. We record the resulting confusion matrix for each round and calculate the precision and recall rates of our classifier.

Having trained the classifier, we measure its ability to generalise by repeating the above experiment, but replacing the test set with the Secondary dataset.

3.5.3.3 Results

After each run of our Random Forest on our Primary dataset, we gather the True Positive (T_P), False Positive (F_P) and False Negative (F_N) rate for each class. We then calculate their precision, defined as $\frac{T_P}{T_P + F_P}$, and recall, defined as $\frac{T_P}{T_P + F_N}$, values. In Fig. 3.7, we see that our average precision and recall across the classes exceeds 0.9 after 10 IATs. Furthermore, after 12 packets our DoS and SQLi data is classified with precision and recall rates of 1.0 and 1.0 and 0.9462 and 0.9322 respectively.

These results do not hold when we test the classifier on our Secondary dataset. As seen in Fig. 3.7, we see a substantial decrease in our average precision and recall rates, achieving a maximum of 0.5923 and 0.5676 respectively. Moreover, after four packets, increasing the number of IATs in our dataset provides little additional benefit. Although some services generalised well, such as IRC-client and IRC-server, others failed to be classified, with every single SMTP feature being classified as HTTP-client. We also see a substantial drop-off in the classification of malicious traffic, with the precision rates of DoS and SQLi data not exceeding 0.6.

These diverging results demonstrate the necessity of dynamic dataset generation for evaluation purposes. Researchers evaluating their methods only on a dataset with fixed properties such as the Primary dataset might receive overoptimistic results. The capability of generating two or more datasets with the same traffic classes, but otherwise differing properties, provides a more realistic evaluation.

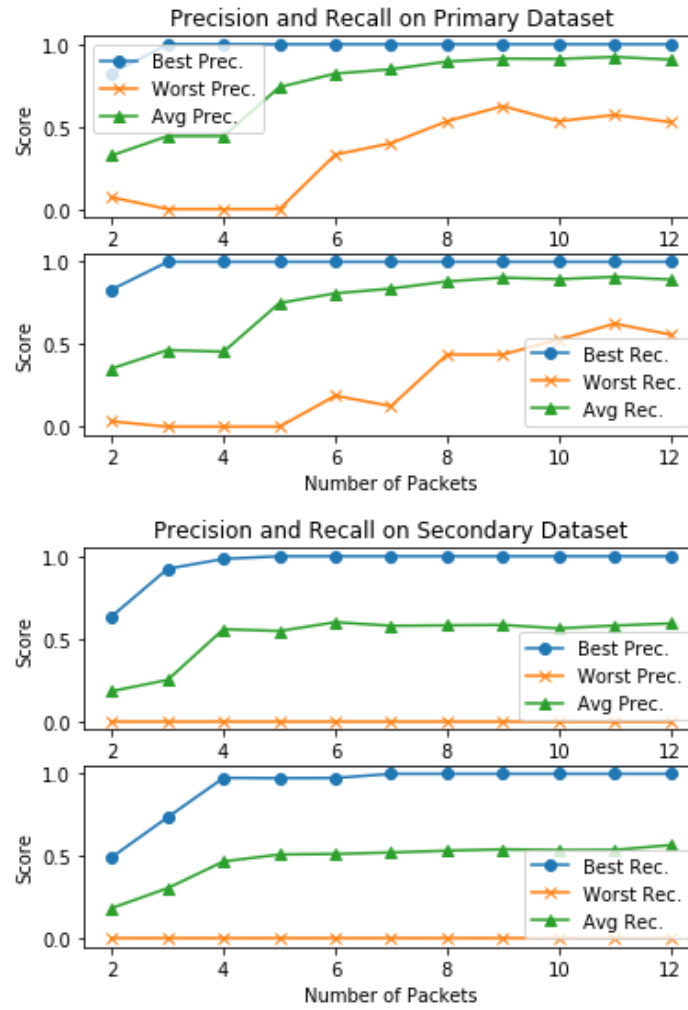


Figure 3.7: Results of Random Forest Classification on Primary dataset (Above) and Secondary dataset (Below)

3.6 Conclusions

In this chapter, we outlined four requirements a modern dataset has to fulfil to strengthen the training of intrusion detection systems. We then proposed a Docker framework capable of generating network intrusion datasets that satisfy these conditions. The major design advantage of this framework are the isolation of traffic scenarios into separate container arrangements, which allows the extension of new scenarios and detailed implementation of subscenarios as well as the capture of ground truth of the computational origins of individual traffic events. Furthermore, containerisation enables the generation of traffic data at scale due to containers being light-weight and easily clonable.

We verified the realism of the generated traffic and the corresponding ground truth

information with two experiments, and demonstrated the usefulness of the framework in another experiment. Presently, our framework consists of 29 scenarios capable of producing benign and malicious network traffic. Several of these scenarios, such as the *BitTorrent* or the *Stepping-Stone* scenario, provide novel traffic data of protocols or behaviours that has not been widely available to researchers previously.

3.6.1 Difficulties and limitations

Our framework is building network traffic datasets from a small-scale level up by coalescing traffic from different fine-grained scenarios together. While this provides great insight into small-scale traffic structures, our framework will not replicate realistic network-wide temporal structures, such as port usage distributions or long-term temporal activity. These quantities would have to be statistically estimated from other real-world traffic beforehand to allow our framework to emulate such behaviour reliably. Other datasets such as UGR-16 use this approach to fuse real-world and synthetic traffic and are currently better suited to build models of large-scale traffic structures.

Working with Docker containers can sometimes complicate the implementation of individual scenarios compared to working with VMs. Although several applications are officially maintained Docker containers that are free from major errors, many do not. For instance, in the *BitTorrent* scenario, most common command line tools, such as `mktorrent`, `ctorrent` and `buildtorrent`, failed to actually produce functioning torrent files from within a container due to Docker's union filesystem. Furthermore, due to the unique way in which we are using these software packages, unusual configuration settings are sometimes needed.

Lastly, capturing `.pcap`-files from each container can quickly exceed available disc space when generating traffic at scale. Depending on specific research requirements, it is advisable to add filtering or feature extraction commands to the scenario execution scripts to enable traffic preprocessing in real-time.

3.6.2 Future work

Although ground truth for particular traffic traces is provided by capturing `.pcap`-files for each container individually, we have not implemented a labelling mechanism yet for the dataset coalescence process. Though not technically difficult, some thought will have to be put how such labels would look like to satisfy different research demands. Furthermore, the Docker platform provides the functionality to collect system logs via

the `syslog` logging driver. We plan on implementing their collection in the future, where they could act either as traffic labels providing more ground truth details, or act as a separate data source that complements the collected traffic.

We paid meticulous attention to enable control over as many traffic impact factors as possible. However, DetGen is currently only offering insufficient control over underlying application-layer implementations such as TLS 1.3 vs 1.2. In theory, it should be unproblematic to provide containers with different implementations, and we are currently investigating how to compile containers in a suitable manner.

While the functionality of DetGen has been verified to some degree in Chapter 4, there are still aspects that need to be properly verified.

Chapter 4

Traffic generation to probe and understand model behaviour

4.1 Introduction

Scientific machine learning model development requires both **model evaluation**, in which the overall predictive quality of a model is assessed to identify the best model, as well as model validation, in which the behaviour and limitations of a model is assessed through targeted **model probing**, as depicted in Fig. 4.1. Model validation is essential to understand how particular data structures are processed, and enables researchers to develop their models accordingly. Data generation tools for rapid model probing such as the *What-If tool* [80] underline the importance of model validation, but are not suitable for providing probing data that resembles the complex structures found in network packet streams.

Machine-learning breakthroughs in many fields have been reliant on a precise understanding of data structure and corresponding descriptive labelling to develop more suitable models. In *automatic speech recognition (ASR)*, tone and emotions can alter the meaning of a sentence significantly. The huge automatically gathered speech datasets however only contain speech snippets and if possible their plain transcripts. While modern speech models are in principle able to learn implicit structures such as emotions without explicit labels, it is impossible to determine the cause for systematic error when they are not. Datasets that contain labelled specialised speech characteristics such as the Ryerson Database of Emotional Speech and Song (RAVDESS) [81] not only allow researchers to identify if their model is susceptible to structural misclassification through targeted probing, but also inspire new methods to capture and

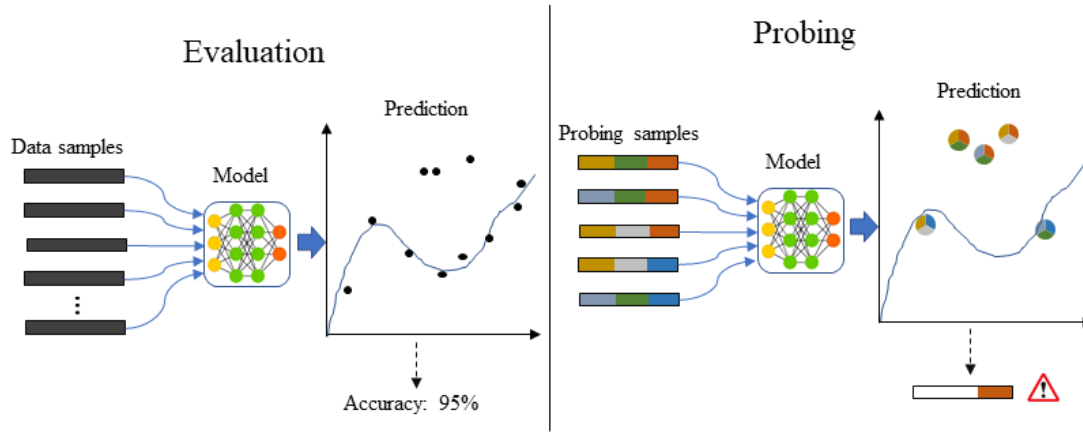


Figure 4.1: Model evaluation and model probing with controlled data characteristics.

understand these implicit structures [82], which in turn leads to design improvements of general speech recognition models [83].

In contrast to ASR, no effort has been made so far to monitor or control the effect of similar factors on network traffic to probe these models for specific microstructures. The current quasi-benchmark NID-datasets such as CICIDS-17, UGR-16 or UNSW-15 pay more attention to the inclusion of specific attacks, protocols, and topologies rather than the documentation of the generated traffic. This situation has so far led researchers to often simply evaluate a variety of ML-models on these datasets in the hope of edging out competitors, without understanding model flaws and corresponding data structures through targeted probing.

In this chapter, we demonstrate how to produce traffic effectively to probe a state-of-the-art traffic classifier, and why a certain degree of generative **determinism** is required for this to isolate the influence of traffic microstructures. The model insights and corresponding performance improvements achieved through probing motivate our experimental examination of various influence factors over microstructures and how to control them during traffic generation.

Thesis context: This chapter builds up on Chapter 3 by extending the examination of DetGen. Here, we examine DetGen’s near-deterministic control over microstructure-shaping factors such as conducted activity, communication failures or network congestion to generate reproducible traffic samples along with corresponding ground-truth labels.

The majority of work in this chapter was published in “Examining traffic microstructures to improve model development” (Henry Clausen and David Aspinall, 2021 [2]) and “Controlling network traffic microstructures for machine-learning model

probing” (Henry Clausen, David Aspinall, and Robert Flood, 2021 [3]).

This chapter discusses the following results:

1. We demonstrate why model probing with controllable traffic microstructure is a crucial step to understand and ultimately improve model behaviour by probing a state-of-the-art LSTM traffic classifier and lowering its false positives five-fold.
2. We examine experimentally how different factors affect traffic microstructures, and how well they can be controlled in a more effective manner when compared to common VM-based traffic generation setups.
3. We examine how DetGen provides accurate control and labels over traffic microstructures, and experimentally demonstrate the level of provided generative determinism to traditional generation set-ups.
4. We discuss requirements for traffic data suitable to probe models pre-trained on a given NIDS-dataset, and demonstrate how to generate probing traffic effectively through DetGen-IDS, a dedicated probing dataset.

4.1.1 Existing datasets and ground-truth information

Real-world NID-datasets such as those from the Los Alamos National Laboratory [39] (LANL) or the University of Grenada [41] provide large amounts of data from a particular network in the form of flow records. Due to the lack of monitoring and traffic anonymisation, it is impossible for researchers to extract detailed information about the specific computational activity associated with a particular traffic sample. Synthetic NID-datasets such as the CICIDS-17 and 18 [84] or the UNSW-NB-15 [85] aim to provide traffic from a wide range of attacks as well as an enterprise-like topology in the form of pcap-files and flow-statistics. While some effort is put in the generation of benign activities using activity scripting or traffic generators, we have seen no attention being spent at monitoring these activities accordingly, which leaves researchers with the limited information available through packet inspection. Furthermore, synthetic datasets can be criticised for their limited activity range, such as the CICIDS-17 dataset where more than 99% of FTP-transfers consist downloading the Wikipedia page for ‘Encryption’ [86], which leads to insufficient structural nuances to for effective training or probing.

4.1.2 Outline

The remainder of the chapter is organised as follows. Section 4.2 and 4.3 motivate the need for probing data with sufficient microstructure control by examining the probing and corresponding improvement of two state-of-the-art intrusion detection models as a motivating example. Section 4.4 proceeds to examine over which traffic characteristics DetGen exerts control and the corresponding control level. Section 4.5 provides details over the design paradigm of DetGen and the resulting advantages over traditional setups, while Section 4.6 discusses the level of control DetGen provides when compared to traditional setups. Section 4.7 discusses how to generate probing data appropriately for pretrained models, and provides an overview over the DetGen-IDS data. Section 4.8 concludes our work.

4.2 Methodology and example

Assume the following problem: You are designing a packet-level traffic classifier which is generating a significant number of false positives, something that is still a common problem for the state-of-the-art [13]. The false positives turn out to be caused by a particular characteristic such as unsuccessful logins or frequent connection restarts. However, existing real-world or synthetic datasets do not contain the necessary information to associate traffic events with these characteristics, which prevents you from identifying the misclassification cause effectively. To address this problem, we need a way to controllably generate and label traffic microstructures caused by these characteristics.

To provide an example, we look at a *Long-Short-Term Memory* (LSTM) network, a deep learning design for sequential data, by Hwang et al. [87], which is designed to classify attacks in web traffic and has achieved some of the highest detection rates of packet-based classifiers in a recent survey [88]. Through probing we will learn that retransmissions in a packet sequence dramatically deplete the model's classification accuracy. We take the following steps:

Step 1: Determine model performance and feed it suitable probing traffic.

Step 2: Examine the correlation between traffic misclassification scores and the generated traffic microstructure labels to find a likely cause.

Step 3: Examine at which latency levels specific connections are misclassified.

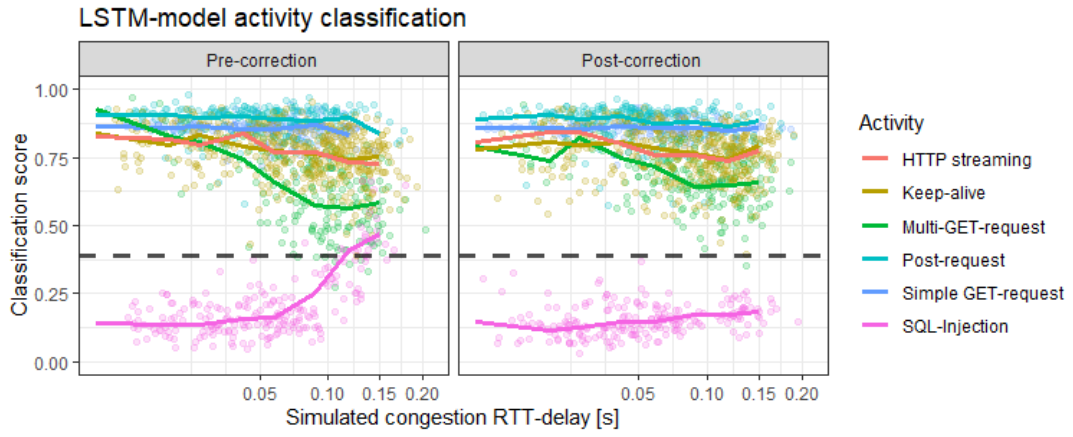


Figure 4.2: Scores for the LSTM-traffic model before and after the model correction.

Step 4: Generate two similar connections, with one exposed to strong packet latency.

Step 5: Show that by removing retransmission sequences in the pre-processing, misclassification is significantly reduced.

Step 1: To detect SQL injections, we train the model on the CICIDS-17 dataset [84] (85% of connections). For the evaluation, we also include a set of HTTP-activities generated by DetGen (7.5%) that mirror the characteristics in the training data, as explained in Section 4.7. In total, we use 30,000 connections for training and for evaluating the model, or slightly under 2 million packets. The initially trained model performs relatively well, with an *Area under curve* (AUC)-score of **0.981**, or a detection/false positive rate¹ of **96%** and **2.7%**. However, to enable operational deployment the false positive rate would need to be several magnitudes lower [89].

Step 2: Now suppose we want to improve these rates to both detect more SQL-injections and retain a lower false positive rate. To start, we explore which type of connections are misclassified most often. We retrieve the classification scores for all connections and measure their linear correlation to the microstructure labels available for the probing data. The highest misclassification ratio was measured for one of the three SQL injection scenarios (19% correlation) and connections with multiple GET-requests (11% correlation). When not distinguishing activities, we measured a high misclassification correlation with simulated packet latency (12%), which we now examine. More details on this exact procedure can be found in (citation currently blinded).

¹tuned for the geometric mean

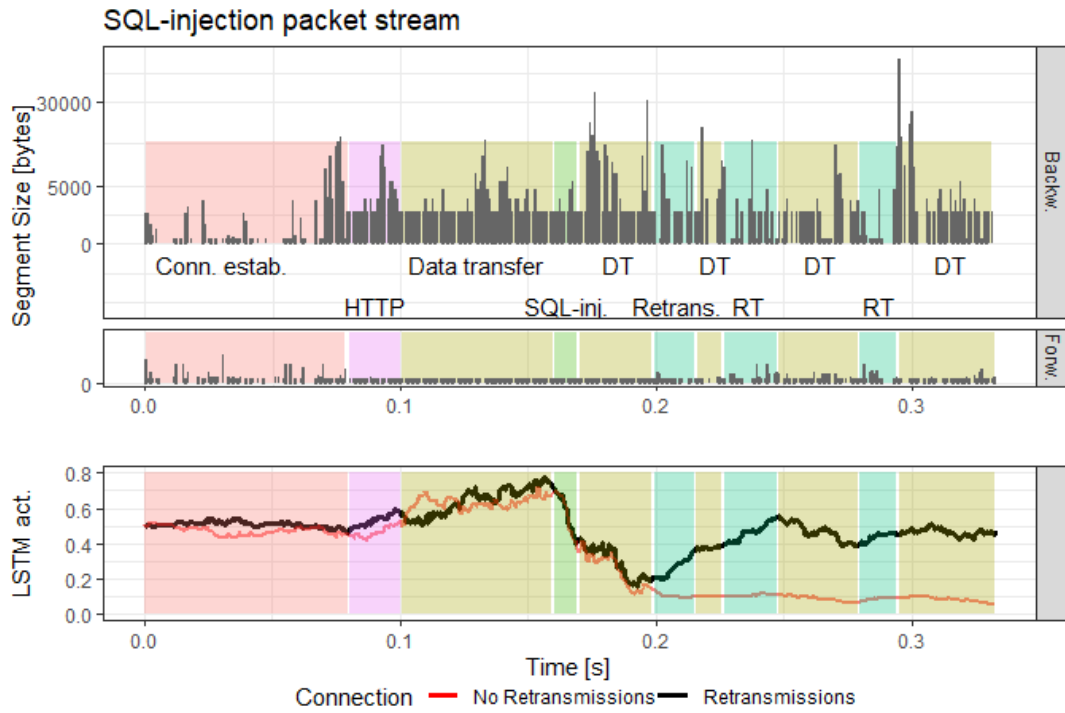


Figure 4.3: LSTM-output activation in dependence of connection phases.

Step 3: Fig. 4.2 depicts classification scores of connections in the probing data in dependence of the emulated network latency. The left panel depicts the scores for the initially trained model, while the right panel depicts scores after the model correction that we introduce further down. The left panel shows that classification scores are well separated for lower congestion, but increased latency in a connection leads to a narrowing of the classification scores, especially for SQL-injection traffic. Since there are no classification scores that reach far in the opposing area, we conclude that congestion simply makes the model lose predictive certainty. Increased latency can both increase variation in observed packet interarrival times (IATs), and lead to packet out-of-order arrivals and corresponding retransmission attempts. Both of these factors can decrease the overall sequential coherence for the model, i.e. that the LSTM-model loses context too quickly either due to increased IAT variation or during retransmission sequences.

Step 4: We use DetGen to generate two similar connections, where one connection is subject to moderate packet latency and corresponding reordering while the other is not. DetGen’s ability to shape traffic in a controlled and deterministic manner allows us to examine the effect of retransmission sequences on the model output and isolate it from other potential influence factors. Fig. 4.3 depicts the evolution of the

LSTM-output layer activation in dependence of difference connection phases for the connection subject to retransmissions. Depicted are packet segment streams and their respective sizes in the forward and backward direction, with different phases in the connection coloured and labelled. Below is the LSTM-output activation while processing the packet streams. The red line shows the output for the connection without retransmissions² as a comparison. Initially the model begins to view the connection as benign when processing regular traffic, until the SQL-injection is performed. The model then quickly adjusts and provides a malicious classification after processing the injection phase and the subsequent data transfer, just as it is supposed to.

The correct output activation is however quickly depleted once the model processes a retransmission phase and is afterwards not able to relate the still ongoing data transfer to the injection phase and return to the correct output activation. When we compare this to the connection without retransmissions, depicted as the red line in Fig. 4.3, we do not encounter this depletion effect. Instead, the negative activation persists after the injection phase.

Step 5: Based on this analysis, we try to correct the existing model with a simple fix by excluding retransmission sequences at the pre-processing stage. This leads to significantly better classification results during network latency, as visible in the right panel of Fig. 4.2. SQL-injection scores are now far-less affected by congestion while scores for benign traffic are also less affected, albeit to a smaller degree. The overall AUC-score for the model improves to **0.997** while tuned detection rates improved to **99.1%** and false positives to **0.345%**, a five-fold improvement from the previous false positive rate of 2.7%.

4.3 Refining the notion of benign traffic for anomaly detection

Next, we show how ground-truth traffic information can help produce more coherent clusters and thus refine the benign traffic model in anomaly-detection. In particular, we will examine a simplified version of *Kitsune* [15], a recent deep learning anomaly-detection model based on stacked autoencoders. *Kitsune*'s AUC-scores surpassed those of other state-of-the-art methods for a variety of attacks, including various types of Botnet traffic and *man-in-the-middle* attacks.

²scaled temporally to the same connection phases

The model takes connection packet streams as input, which are pushed through an artificial information bottleneck before reconstruction, which forces the model to learn and compress reoccurring traffic structures. The compressed connection representation is essentially a positional projection into a lower-dimensional vector space, where spatial boundaries around benign traffic can be drawn. For demonstration purposes, we use a widely-used clustering approach for anomaly-detection rather than *Kitsune*'s more complex ensemble method. Here, anomalous outliers are detected using the Mahalanobis-distance of a projected connection from identified cluster centres. Benign traffic should ideally be distributed evenly around the cluster centres to allow a tight borders and good separation from actual abnormal behaviour.

Unstructured datasets such as the CAIDA traffic traces assumably contain too much abnormal behaviour to train an anomaly-detection model, which is why we train the model on benign traffic from the CICIDS-17 [84] intrusion detection dataset (80%). Again, we add 20% probing traffic consists of HTTP, FTP, SSH, and SMTP communication, using a wide spectrum of settings for examination purposes. Attack data for the evaluation was again provided through the CICIDS-17 dataset, and includes access attacks such as SQL-injections or Brute-Forcing, as well as Mirai botnet traffic. We train the model with in total 150,000 connections.

4.3.1 Projection coherency evaluation

Like many approaches that generate representations of benign traffic for anomaly detection, *Kitsune* projects traffic events into a vector-space where traffic clusters and similarities become more apparent. In order for the projection to accurately capture important traffic structures, this projection should be consistent, i.e. traffic events with similar origins and characteristics should be projected to similar positions rather than be dispersed throughout the vector space [90].

To verify the models projection consistency, we generate traffic from near-identical conditions to provide certainty on the expected traffic similarities. We generate a small dataset that consists of HTTP-requests, file-synchronisation, and Botnet communication. For each of the three traffic types we fix four settings that vary in the performed activity and network latency, with the traffic shaping described in Section 4.4 being held constant within each setting except for small variations in the transmitted message or file. Table 4.1 summarises the traffic for each setting.

We verify if traffic samples within each group are projected to similar areas by

Label	HTTP	File-Sync	Mirai-C&C
1	Get-req. NGINX, low lat.	Two hosts, low lat.	Command 1, low lat.
Results:	0.14 , 0.45	0.19 , 0.27	0.03 , 0.06
2	Multi-req. NGINX, low lat.	Four hosts, low lat.	Command 2, low lat.
Results:	0.32 , 0.45	0.15 , 0.33	0.03 , 0.04
3	Post-req. Apache, high lat.	Two hosts, high lat.	Command 3, high lat.
Results:	0.17 , 0.28	0.16 , 0.28	0.02 , 0.04
4	Multi-req. Apache, high lat.	Four hosts, high lat.	Command 4, high lat.
Results:	0.53 , 2.51	0.71 , 1.31	0.03 , 0.05

Table 4.1: Outline of the traffic settings for examining projection consistency. The numbers below each setting describe the measured Mahalanobis-distances (blue:average, red:maximal) for the corresponding projections.

measuring the average and maximum Mahalanobis-distance to quantify the overall dispersion of the samples. The results are displayed in Table 4.1 and depicted in Fig. 4.4. The first thing to notice is that the model projects samples from each group within the same cluster, thus confirming the capture of a coarse traffic structure. When looking at the traffic dispersion and the corresponding Mahalanobis-distance measurements, we notice that the *multi-request HTTP* traffic as well as the *file-synchronisation* between multiple computers is much further dispersed than in the other settings, especially when exposed to more latency. We also find that the corresponding dimension, x_3 , with the most projected dispersion seems to be the same for each of the four settings. This suggests that the cause for the dispersion is the same for the different traffic types.

We now focus on the influence of input features on the projected positions exclusively in the x_3 -direction. Here, we can again perform a simple correlation analysis between different the input feature values and the corresponding x_3 -value. We observe that the arrival time of packet bears the most correlation (5.4%) for the selected settings. We also see that this influence is concentrated primarily on connections that are opened shortly after a previous connection, with the temporal separation between these

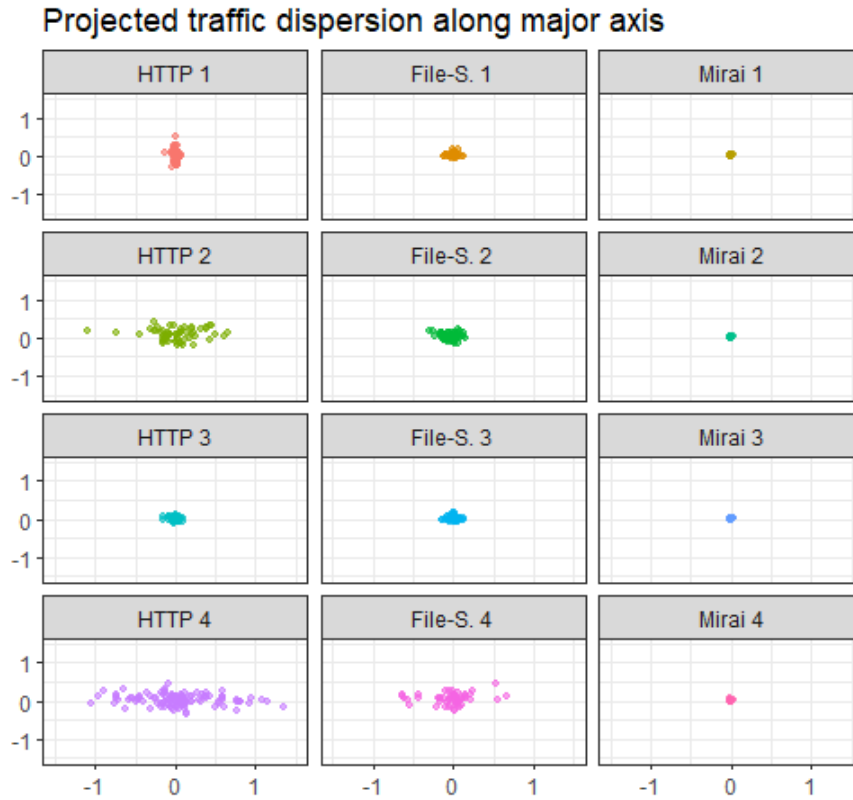


Figure 4.4: Dispersion of projected traffic samples from each setting, plotted along the two most dispersed axes.

two connections apparently being the primary cause for the spread on the x_3 -axis. The connection interarrival times are naturally an important feature for *Kitsune* to detect attacks such as *Man-in-the-Middle*, which could explain the weight this feature plays in the projection process.

4.3.2 Investigating individual cluster incoherences

When examining false-positive and corresponding anomaly scores, we noticed that the model often classifies Brute-Force Web attacks as benign and some HTTP-traffic as anomalous. When examining the projected location of the corresponding connections, we see that most of this HTTP-traffic as well as the Brute-Force attack traffic lie near a particular cluster, depicted in Fig. 4.5. A significant portion of traffic in that cluster seems to be spread significantly more across the cluster axis than the rest of the traffic in that cluster, leading to an inflated radius that partially encompasses Brute-Force traffic.

When cross-examining the traffic in this cluster with the probing data, we see that

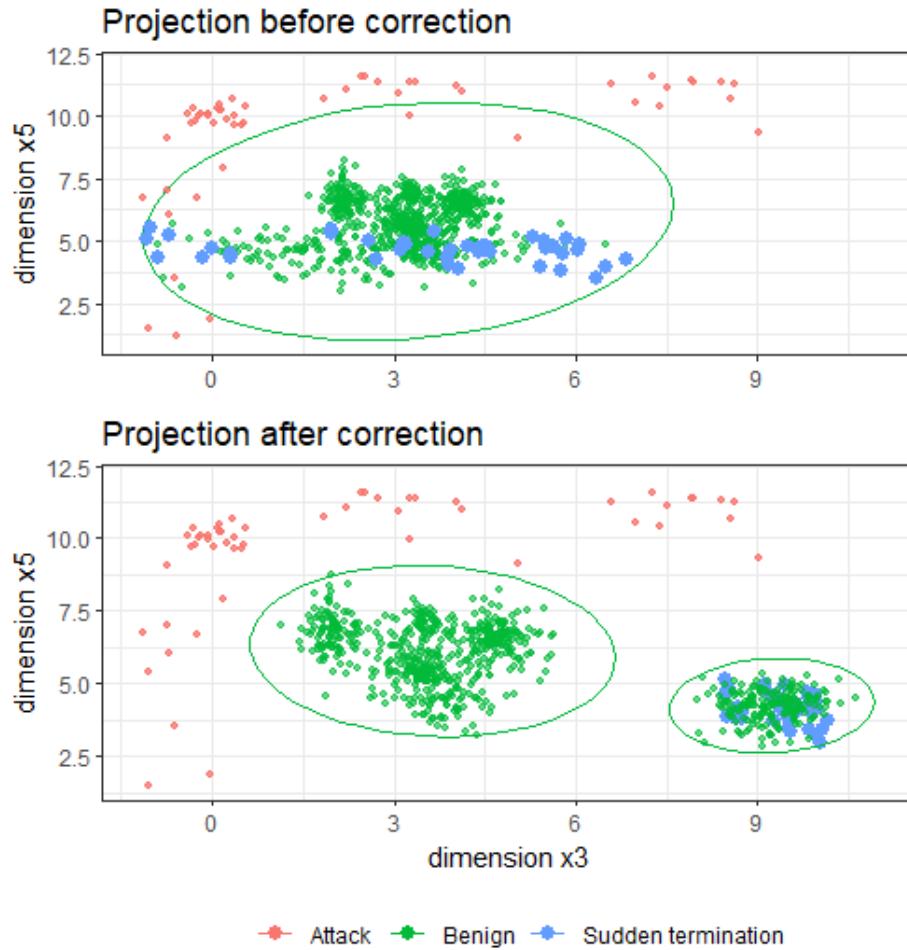


Figure 4.5: Scores for the LSTM-traffic classification model in dependence of simulated network congestion, along with the classification threshold

HTTP-traffic with the label "Sudden termination" are distributed across the cluster axis in a similar fashion, also depicted in Fig. 4.5, suggesting the conclusion that this type of traffic causes the inflated cluster radius. DetGen generates traffic with the label "Sudden termination" as half-open connections which were dropped by the server due to network failure. One defining characteristic of such connections are that they are not closed with a termination handshake using FIN-flags. To better capture this defining characteristics in the modelling process, we included an additional feature attached to the end of a packet sequence that indicates a proper termination with FIN-flags in the modelling process. The newly trained model now projects "Sudden termination" connections into a different cluster, which leads to a far better cluster coherence. The detection rate on Brute-Force attack traffic could thus be improved from **89.7%** to **94.1%**.

4.4 Traffic microstructures and their influence factors

The biggest and most obvious influence on traffic microstructures is the choice of the application layer protocols. For this reason, the range of protocols is often used as a measure for the diversity of a dataset. However, while the attention to microstructures in current NID-datasets stops here, computer communication involves a myriad of other different computational aspects that shape observable traffic microstructures. Here, we highlight and quantify the most dominant ones, which will act as a justification for the design choices we outline in Section 4.5.1. We look at both findings from previous work as well as our own experimental results.

1. Performed task and application. The conducted computational task as well as the corresponding application ultimately drives the communication between computers, and thus hugely influences characteristics such as the direction of data transfer, the duration and packet rate, as well as the number of connections established. These features are correspondingly used extensively in application fingerprinting, such as by Yen et al. [91] or Stober et al. [92].

2. Application layer implementations. Different implementations for TLS, HTTP, etc. can yield different computational performance and can perform handshakes differently and differ in multiplexing channel prioritisation, which can significantly impact IAT times and the overall duration of the transfer, as shown in a study by Marx et al. [93] for the QUIC/HTTP3 protocol³.

3. LAN and WAN congestion. Low available bandwidth, long RTTs, or packet loss can have a significant effect on TCP congestion control mechanisms that influence frame-sizes, IATs, window sizes, and the overall temporal characteristic of the sequence, which in turn can influence detection performance significantly as shown in Section 4.2.

4. Host level load. In a similar manner, other applications exhibiting significant computational load (CPU, memory, I/O) on the host machine can affect the processing speed of incoming and outgoing traffic, which can again alter IATs and the overall duration of a connection. An example of this is visible in Fig. 4.6, where a FTP-client sends significantly fewer PUSH-packets when under heavy computational load. Colours

³Fig. 2 in [93] illustrates these differences in a nice way

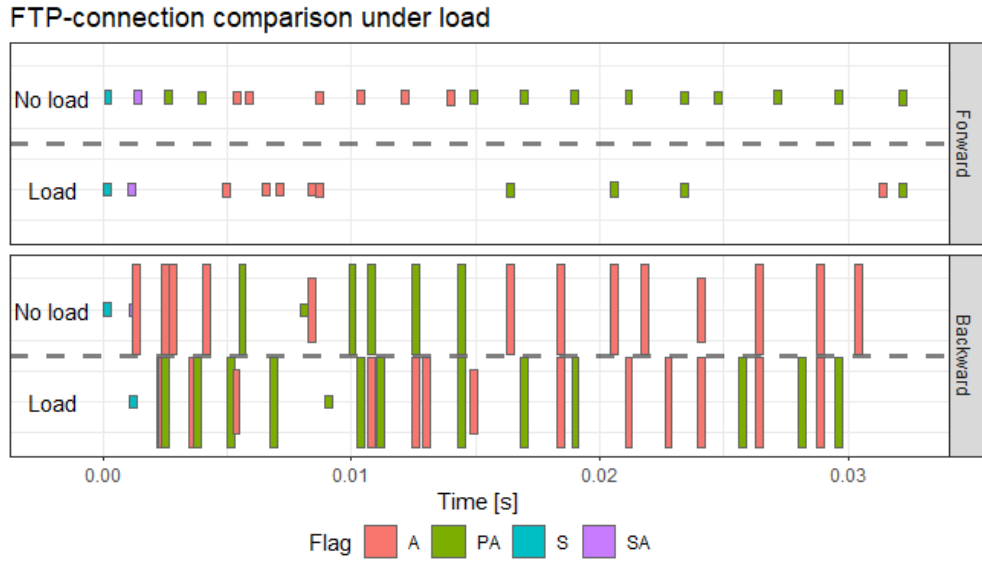


Figure 4.6: Packet-sequence similarity comparison under different load.

indicate packet flags while the height of the packets indicates their size. This effect is dependent on the application layer protocol, where at a load number of 3.5 we see about 60% less upstream data-packets while the downstream is only reduced by 10%, compared to HTTP where both downstream and upstream packet rates are throttled by about 40%.

5. Caching/Repetition effects. Tools like cookies, website caching, DNS caching, known hosts in SSH, etc. remove one or more information retrieval requests from the communication, which can lead to altered packet sequences and less connections being established. For caching, this can result in less than 10% of packets being transferred, as shown by Fricker et al. [94].

6. User and background activities. The choice and usage frequency of an application and task by a user, sometimes called *Pattern-of-Life*, governs the larger-scale temporal characteristic of a traffic capture, but also influences the rate and type of connections observed in a particular time-window [95]. The mixing of different activities in a particular time-window can severely impact detection results of recent sequential connection-models, such as by Radford et al. [96] or by Clausen et al. [4]. To quantify this effect, we look at FTP-traffic in the CICIDS-17 dataset. As explained in Section 4.1.1, the FTP-traffic overwhelmingly corresponds to the exact same isolated task, and should therefore spawn the same number of connections in a particular time window. However, we observe additional connections from other activities within a 5-second

window for 68% of all FTP-connections, such as depicted in Table 4.2, which contains FTP-, HTTPS- and DNS-, as well as additional unknown activity.

Time	Source-IP	Destination-IP	Dest. Port
13:45:56.8	192.168.10.9	192.168.10.50	21
13:45:56.9	192.168.10.9	192.168.10.50	10602
13:45:57.5	192.168.10.9	69.168.97.166	443
13:45:59.1	192.168.10.9	192.168.10.3	53
13:46:00.1	192.168.10.9	205.174.165.73	8080

Table 4.2: 5-second window for host 192.168.10.9 in the CICIDS-17 dataset.

Other prominent factors that we found had less effect on traffic microstructures include:

7. Networking stack load. TCP or IP queue filling of the kernel networking stack can increase packet waiting times and therefore IATs of the traffic trace, as shown by [97]. In practice, this effect seems to be constrained to large WAN-servers and routers. When varying the stack load in otherwise constant settings on an Ubuntu-host, we did not find any notable effect on packet sequences when comparing the corresponding traffic with a set of three similarity metrics. More details on this setting and the metrics can be found in Section 4.6.

8. Network configurations. Network settings such as the MTU or the enabling of TCP Segment Reassembly Offloading have effects on the captured packet sizes, and have been exploited in IP fragmentation attacks. However, these settings have been standardised for most networks, as documented in the CAIDA traffic traces [43].

We designed DetGen to control and monitor factors 1-6 to let researchers explore their impact on their traffic models, while omitting factors 7 and 8 for the stated reasons.

4.5 How DetGen generates precisely controlled data

As we will demonstrate in Section 4.6, containers provide significantly more isolation of programs from external effects than regular OS-level execution. This isolation

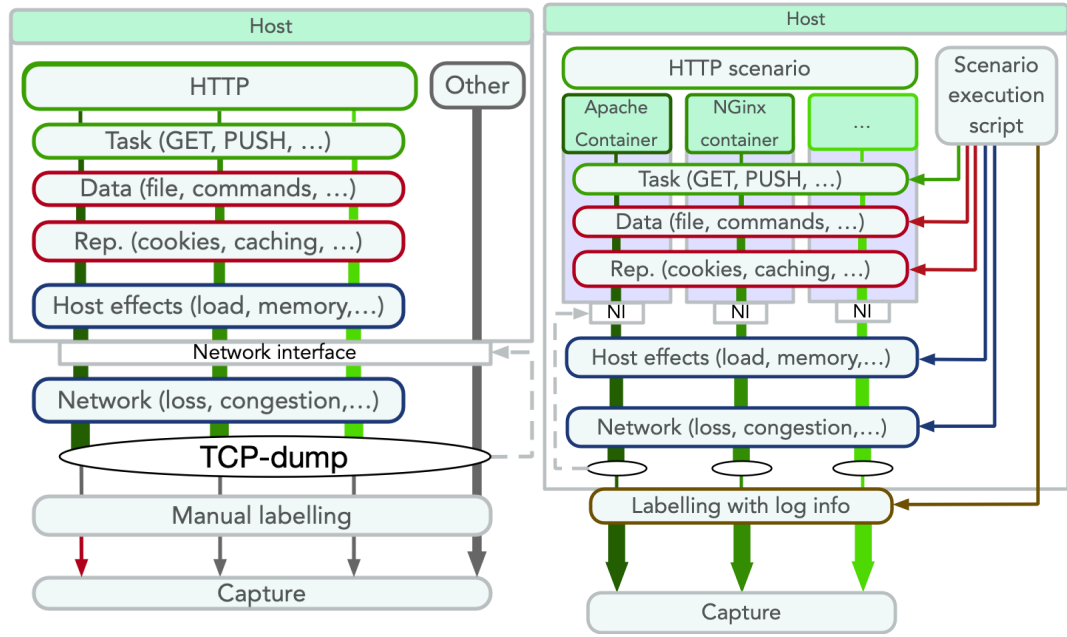


Figure 4.7: Traditional traffic-generation-setups (left), and DetGen (right).

enables us to monitor processes better and create more accurate links between traffic events and individual activities than on a virtual machine where multiple processes run in parallel and generate traffic. The corresponding one-to-one correlation between processes and network traces allows us to capture traffic directly from the process and produce labelled datasets with granular ground truth information.

Additionally, containers are specified in an image-layer, which is unaffected during the container execution. This allows containers to be run repeatedly whilst always starting from an identical state, allowing a certain level of **determinism** and reproducibility in the data generation.

4.5.1 Simulation of external influence

4.5.1.1 Caching/Cookies/Known server.

Since we always launch containers from the same state, we prevent traffic impact from **repetition effects** such as caching or known hosts. If an application provides caching possibilities, we implement this as an option to be specified before the traffic generation process.

4.5.1.2 Network effects.

Communication between containers takes place over a virtual bridge network, which provides far higher and more reliable throughput than in real-world networks [68]. To retard and control the network reliability and congestion to a realistic level, we rely on *NetEm*, an enhancement of the Linux traffic control facilities for emulating properties of wide area networks from a selected network interface [64].

We apply *NetEm* to the network interface of a given container, providing us with the flexibility to set each container's network settings uniquely. In particular, packet delays are drawn from a Paretonormal-distribution while packet loss and corruption are drawn from a binomial distribution, which has been found to emulate real-world settings well [98]. Distribution parameters such as mean or correlation as well as available bandwidth can either be manually specified or drawn randomly before the traffic generation process.

4.5.1.3 Host load.

We simulate excessive computational load on the host with the tool *stress-ng*, a Linux workload generator. Currently, we only stress the CPU of the host, which is controlled by the number of workers spawned. Future work will also include stressing the memory of a system. We have investigated how stress on the network sockets affects the traffic we capture without any visible effect, which is why we omit this variable here.

4.6 Verifying the generative determinism of DetGen

We now assess the claim that DetGen controls the outlined traffic influence factors sufficiently, and how similar traffic generated with the same settings looks like. We also demonstrate that this level of control is not achievable on regular VM-based NIDS-traffic-generation setup.

To do so, we generate traffic from three traffic types, namely HTTP, file-synchronisation, and botnet-C&C, each in four configurations that varied in terms of conducted activity, data/credentials as well as the applied load and congestion. Within each configuration all controllable factors are held constant to test the experimental determinism and reproducibility of DetGen's generative abilities. As a comparison, we use a regular VM-based setup, where applications are hosted directly on two VMs that communicate over a virtual network bridge that is subject to the same *NetEm* effects as DetGen,

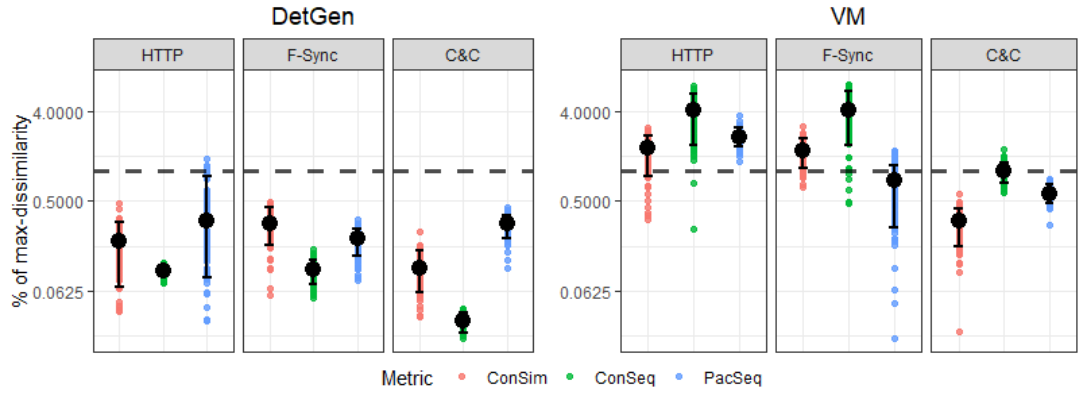


Figure 4.8: Dissimilarity scores for DetGen and a regular VM-setup, on a log-scale.

such as depicted in Fig. 4.7. Such a setup is for example used in the generation of the UGR-16 data [41].

To measure how similar two traffic samples are, we devise a set of similarity metrics that measure dissimilarity of overall connection characteristics, connection sequence characteristics, and packet sequence characteristics:

- **Overall connection similarity** We use the 82 flow summary statistics (IAT and packet size, TCP window sizes, flag occurrences, burst and idle periods) provided by CICFlowMeter [99], and measure the cosine similarity between connections, which is also used in general traffic classification [100].
- **Connection sequence similarity** To quantify the similarity of a sequence of connections in a retrieval window, we use the following features to describe the window, as used by Yen et al. [91] for application classification: The number of connections, average and max/min flow duration and size, number of distinct IP and ports addresses contacted. We then again measure the cosine similarity based on these features between different windows.
- **Packet sequence similarity** To quantify the similarity of packet sequences in handshakes etc., we use a Markovian probability matrix for packet flags, IATs, sizes, and direction conditional on the previous packet. We do this for sequences of 15 packets and use the average sequence likelihood as this accommodates better for marginal shifts in the sequence.

We normalise all dissimilarity scores by dividing them by the maximum dissimilarity score measured for each traffic type to put the scores into context. For each configuration, we generate 100 traffic samples and apply the described dissimilarity

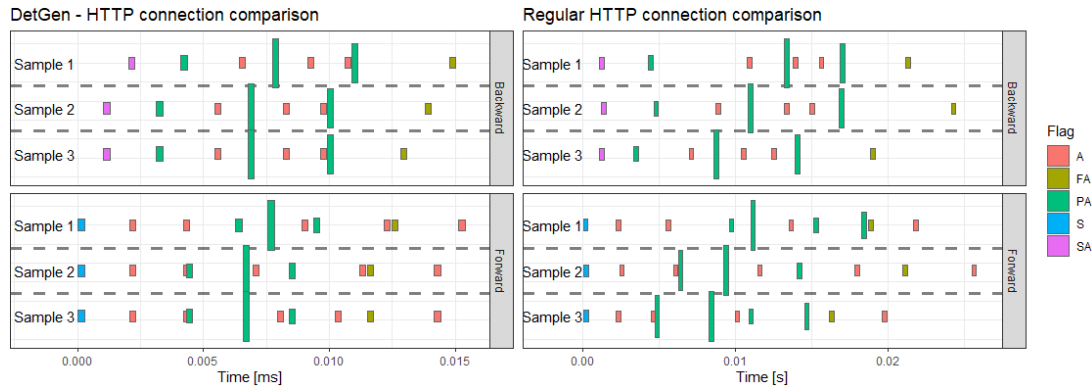


Figure 4.9: Packet-sequence similarity comparison for HTTP-activity for DetGen and a regular setting.

measures to 100 randomly drawn sample pairs. Fig. 4.8 depicts the resulting dissimilarity scores on a log-scale.

The DetGen-scores yield consistently less than 1% of the dissimilarity observed on average for each activity. Scores are especially low when compared to traffic groups collected in the VM setting, which are consistently more dissimilar, in particular for connection-sequence metrics, where the average dissimilarity is more than 30 times higher than for the DetGen setting. Manual inspection of the VM-capture showed that high dissimilarity is caused by additional flow events from background activity (OS and application HTTP, NTP, DNS, device discovery) being present in about 24% of all captures. . While sequential dissimilarity is roughly the same for the DetGen- and the VM-settings, overall connection similarity for the VM-setting sees significantly more spread in the dissimilarity scores when computational load is introduced.

Fig. 4.9 depicts an exemplary comparison between HTTP-samples generated using DetGen versus generation using the VM-setup. Colours indicate packet flags while the height of the packets indicates their size. Even though samples from DetGen are not perfectly similar, packets from the VM-setup are subject to more timing perturbations and reordering as well as containing additional packets. Additionally, the packet sizes vary more in the regular setting.

These results confirm that DetGen exerts a high level of control over traffic shaping factors while providing sufficient determinism to guarantee ground-truth traffic information.

4.7 Reconstructing an IDS-dataset for efficient probing

Moving towards a more general dataset constructed to apply this probing methodology, we constructed *DetGen-IDS*. This dataset is suitable to quickly probe ML-model behaviour that were trained on the CICIDS-17 dataset [84]. The dataset mirrors properties of the CICIDS-17 data to allow pre-trained models to be probed without re-training. The *DetGen-IDS* data therefore serves as complementary probing data that provides microstructure labels and a sufficient and controlled diversity of several traffic characteristics that is not found in the CICIDS-17 data.

We focus on mirroring the following properties from the CICIDS-17 data:

1. **Application layer protocols (ALP)**
2. **ALP implementations**
3. **Typical data volume for specific ALPs**
4. **Conducted attack types**

Extracting more information on characteristics such as conducted activities of current NID-datasets is difficult for the reasons explained in 4.1.1. However, our examination shows that aligning these high-level features with the original training data helps to significantly reduce the validation error of a model on the probing data.

We then took the following steps to extract the necessary information from the CICIDS-17 data and implement the traffic-generation process accordingly:

1. The primary ALPs in the dataset can be identified using their corresponding network ports. We ordered connections by the frequency of their respective port, and excluded connections that do not transmit more than 15 packets per connection as these do not provide enough structure to create probing data from it. This leaves us with the ALPs *HTTP/SSL*, *SMTP*, *FTP*, *SSH*, *SQL*, *SMB*, and *NTP*. We had already implemented traffic scenarios for each of them except SMB and LDAP, which we then added to the catalogue described in Section 3.4.1. Table 4.3 displays the frequency of the most common ALPs in the CICIDS-17 along with their average size and packet number per connection and how we adopted them in the *DetGen-IDS* data.

2. Most of the used ALP implementations, such as *Apache* and *Ubuntu Webserver* for HTTP, could be gathered from the description of the CICIDS-17 dataset. When this was not the case, it is mostly possible to gather this information by inspecting a few negotiation packets for the corresponding ALP with Wireshark to identify the

TLS version or the *OpenSSH*-client. The correct ALP implementation can then be included in the traffic generation process by simply identifying and including a Docker-container that matches the requirements, which is explained more in Section 3.4.2.

ALP	Port	CICIDS-17			DetGen-IDS	
		Av. Conn. Size	Av. Packets /Conn.	# Packets	# Packets	# Activities
HTTP	80	131626.4	120.4	26631853	724032	7
HTTPS	443	24637.5	36.7	18531661	432104	7
DNS	53	286.2	3.6	3515510	-	-
SSH	22	4699.6	40.9	430380	379421	13
LDAP	389	5429.2	22.3	133471	94587	3
FTP	21	311.3	41.7	121472	183587	9
NetBIOS	137	773.6	14.3	111341	-	-
SMB	445	12941.5	61.9	88175	47945	3
NTP	123	157.0	3.2	73057	1243	1
SMTP	465	2663.5	21.5	77650	104967	3
Kerberos	88	2687.7	6.9	38262	-	-
mDNS	5353	3685.5	35.5	24592	-	-

Table 4.3: Common ALPs in CICIDS-17 data

3. Since the total size of a connection is one of the most significant features for its classification, we restrict connections in the DetGen-IDS data to cover the same range as their counterparts in the CICIDS-17 data. For this, we extracted the maximum and minimum connection size for each ALP in the benign data and use it as a cut-off to remove all connections from the DetGen-IDS data that do not meet this requirement.

4. Included attacks are well documented in the CICIDS-17 description. These include *SQL-injections*, *SSH-brute-force*, *XSS*, *Botnet*, *Heartbleed*, *GoldenEye*, and *SlowLoris*. We aim to cover as many of these attack types in the DetGen-IDS data as well as adding them to the overall DetGen-attack-catalogue. We were not able to cover all attacks though as DetGen either did not provide the necessary network topology to conduct the attack, such as for port-scanning, or the attack types are not implemented in the catalogue of scenarios yet.

In addition to the pcap-files, we used the *CICFlowMeter* to generate the same 83 flow-features as included in the CICIDS-17 data. Table 4.3 displays the content and statistics of the DetGen-IDS data.

In Fig. 4.10, we compare the validation error of a recent LSTM-model for network intrusion detection by Clausen et al. [4] on the DetGen-IDS data to demonstrate that a model trained on the CICIDS-17 data is able to perform well without retraining. We distinguish models when trained exclusively on the CICIDS data (green), and when

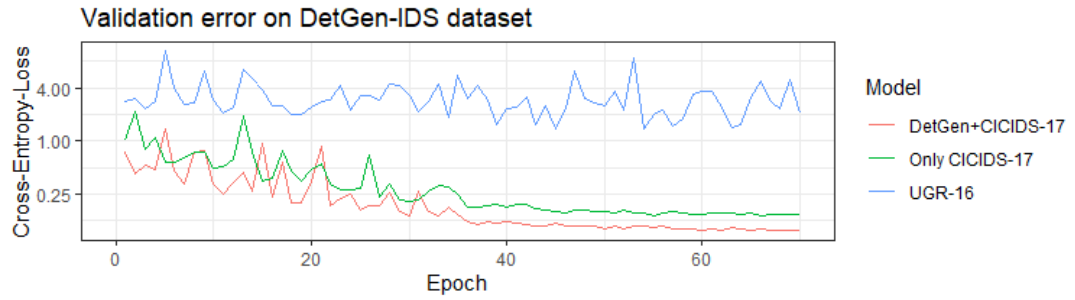


Figure 4.10: Validation errors of LSTM-model [4] on DetGen-IDS data.

also trained on the probing data (red). Even though the validation error is slightly higher when only trained on the CICIDS data, the difference is almost negligible compared to the error resulting from a model trained on a completely different dataset (UGR-16 [41], blue). This does not fully prove that every model is able to transfer observed structures between the two datasets, but it gives an indicator that they mirror characteristics.

4.8 Conclusions

In this chapter, we examined how the generation of traffic with controllable and extensively labelled traffic microstructures can aid the probing of machine-learning-based traffic models. For this, we demonstrated the impact that probing with carefully crafted traffic microstructures can have for improving a model with a state-of-the-art LSTM-traffic-classifier with a detection rate that improved by more than 3% after understanding how the model processes excessive network congestion.

To verify DetGen’s ability to control and monitor traffic microstructures, we performed experiments in which we quantified the experimental determinism of DetGen and compared it to traditional VM-based capture setups. Our similarity metrics indicate that traffic generated by DetGen is on average 10 times, and for connection sequences up to 30 times more consistent.

We are also releasing DetGen-IDS, a substantial dataset suitable for probing models trained on the CICIDS-17 dataset. This data should make it easier for researchers to understand where their model fails and what traffic characteristics are responsible to subsequently improve their design accordingly.

4.8.1 Difficulties and limitations:

While the control of traffic microstructures helps to understand packet- or connection-level models, it does not replicate realistic network-wide temporal structures. Other datasets such as UGR-16 [41] or LANL-15 [40] are currently better suited to examine models of large-scale traffic structures.

While controlling traffic shaping factors artificially helps at identifying the limits and weak points of a model, it can exaggerate some characteristics in unrealistic ways and thus alter the actual detection performance of a model.

The artificial randomisation of traffic shaping factors can currently not completely generate real-world traffic diversity. This problem is however more pronounced in commonly used synthetic datasets such as CICIDS-17, where for example most FTP-transfers consist of a client downloading the same text file.

Discussions about the implications of the model correction proposed in Section 4.2 are above the scope of this chapter, and there likely exist more complex and suitable solutions.

Chapter 5

CBAM: An anomaly detection model for traffic microstructures

5.1 Introduction

Remote access attacks are used to gain control or access information on remote devices by exploiting vulnerabilities in network services, and are involved in many of today's data breaches [101]. A recent survey [13] showed that these attacks are detected at significantly lower rates than more high-volume probing or DoS attacks. To address this, we present **CBAM**, a short-term contextual **bidirectional anomaly model** of network flows, which improves detection rates of remote access attacks significantly by learning reoccurring flow sequence microstructures. The underlying idea of CBAM is to capture probability distributions over sequences of network flows that quantify their overall likelihood, much like a language model. CBAM is based on deep bidirectional LSTM networks.

Recently, deep learning models such as LSTMs have been a popular tool in network intrusion detection [102, 103, 96]. However, persistent failings in evaluations have made it difficult to assess the performance and real-world applicability of currently proposed methods to access attack detection, and have lead to a chaotic and convoluted NIDS landscape [13].

To avoid these pitfalls and demonstrate that our approach delivers a significant improvement in detection rates and real-world applicability, we evaluated our model carefully on three modern network intrusion detection datasets. Furthermore, we reimplemented and evaluated three state-of-the-art methods on these datasets and compared their performance against ours.

We carefully select the input parameters for our model based on the interdependence we observe between them in exemplary flow sequences, and increase the model complexity in terms of depth and input embedding compared to preceding models. This enables us to detect remote access attacks at a false positive rate of 0.16%, a rate at which none of the comparison models are able to detect any attacks reliably.

We also discuss specific design choices and how they enable effective modelling of specific traffic microstructures to boost performance. The evaluation of existing deep learning models so far has generally been agnostic to particular microstructures of the modelled traffic and fails to explain where and why the corresponding model fails to classify traffic properly. In this chapter, we therefore apply the similar model probing techniques as demonstrated in Chapter 4 to validate the undertaken design steps.

Thesis context: This chapter turns from the generation and examination of traffic microstructures to the design of an anomaly detection model that leverages specific microstructures for detection. It also examines which conclusions can be made on the distribution and similarity of microstructures from the model output.

This chapter largely consists of work published in “Better Anomaly Detection for Access Attacks Using Deep Bidirectional LSTMs” (H. Clausen, G. Grov, M. Sabate, and D. Aspinall, 2020 [4]) and “CBAM: A Contextual Model for Network Anomaly Detection” (H. Clausen, G. Grov, and D. Aspinall, 2021 [5]).

5.1.1 Outline

The remainder of this chapter is organised as follows: Section 5.2 provides a motivation for short-term flow models and their benefits for the detection of access attacks. Section 5.3 provides an overview of common evaluation pitfalls that prevent a fair comparison of corresponding NID models. Section 5.4 explains the methodology and architecture of CBAM as well as the data preprocessing. Section 5.5 describes the problems with network traffic datasets which previous methods were evaluated on, and explains the advantages of our selection of datasets. We also describe how and why traffic from particular hosts is selected, and how training and test data are constructed. Section 5.6 discusses our detection rates on attack traffic in the CICIDS-17 and LANL-15 data, and examines how and why CBAM is able to identify these attacks. Section 5.7 discusses the false positive rate on benign traffic, and examines long-term stability of flow structures and corresponding model performance. It also provides details on the score distributions and the type of traffic that is both predicted

accurately and inaccurately, as well as the influence of the training data size on these predictions. Section 5.8 discusses the reason and measured benefit of specific design steps that increase model complexity. Section 5.11 concludes our results.

5.2 Overview

Src	Dst	DPort	bytes	# packets	Src	Dst	DPort	bytes	# packets
A	B	80	247956	315	D	C	N33	600	5
A	B	80	7544	13	C	D	445	77934	1482
A	B	80	328	6	D	C	N33	600	5
A	B	80	2601	10	C	D	445	5202	10
A	B	80	328	6	(b) Benign SMB, C=C6267, D=C754				
A	B	80	328	6	Src	Dst	DPort	bytes	# packets
A	B	80	380	7	C	D	445	4106275	2830
A	B	80	328	6	C	D	445	358305611	242847
⋮									
(a) XSS-attack, A=192.168.10.50, B=172.16.0.1					(c) <i>Pass-the-hash</i> attack via SMB				

Table 5.1: The left side depicts a flow sequence from an XSS-attack. The right side depicts a benign SMB-sequence (top), and a sequence from a *Pass-the-hash* attack via the same SMB service.

In verbal or written speech, we expect the words “I will arrive by ...” to be followed by a word from a smaller set such as “car” or “bike” or “5pm”. Similarly, on an average machine we may expect DNS lookups to be followed by outgoing HTTP/HTTPS connections. These short-term structures in network traffic are a reflection of the computational order of information exchange. Attacks that exploit vulnerabilities in network communication protocols often achieve their target by deviating from the regular computational exchange of a service, which should be reflected in the generated network pattern.

Table 1(a) depicts a flow sequence from an XSS-attack. Initial larger flows are followed by a long sequence of very small flows which are likely generated by the embedded attack script trying to download multiple inaccessible locations. Flows of this

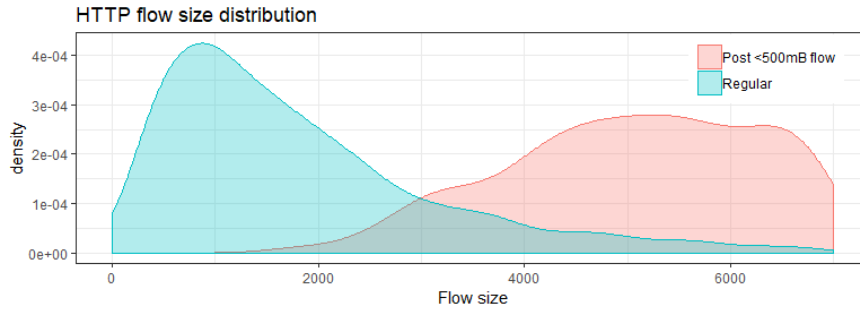


Figure 5.1: HTTP flow size distribution overall, and if preceded by an HTTP flow smaller than 500 bytes.

size are normally immediately followed by larger flows, as depicted in Fig. 5.1, which makes the repeated occurrence of small HTTP flows in this sequence very unusual.

Table 1(b) depicts a regular SMB service sequence while Table 1(c) depicts a *Pass-the-hash* attack via the same SMB service. As shown, the flows to port N33 necessary to trigger the communication on the SMB port are missing while the second flow is significantly larger than any regular SMB flows due to it being misused for exfiltration purposes.

The underlying idea of CBAM is to predict probabilities of connections in a host’s traffic stream conditional on adjacent connections. The probabilities are assigned based on the connection’s protocol, network port, direction, and size, and the model is trained to maximise the overall predicted probabilities.

To assign probabilities, we map each connection event to two discrete sets of states, called vocabularies, according to the protocol, the network port, and the direction of the connection for the first, and according to number of transmitted bytes for the second. The size of the vocabulary is chosen large enough to capture meaningful structures without capturing rare events that can deteriorate prediction quality. We feed these vocabularies into a deep bidirectional LSTM (long short-term memory) network that takes bivariate sequences of mapped events as input to efficiently capture the conditional probabilities for each event.

CBAM acts as an anomaly-detection model that learns short-term microstructures in benign flow sequences and identifies malicious sequences as deviations from these structures. By predicting probabilities of flows in benign flow sequences, CBAM is trained in a self-supervised way on strictly benign traffic. In contrast to classification-based training, CBAM does not require labelled attack traffic in the training data and is thus not affected by typical class imbalances in network intrusion datasets.

5.3 Evaluation pitfalls

According to Nisioti et al. [13], the trustworthiness of published low volume access attack detection rates is debatable due to evaluation shortcomings. We designed our evaluation to avoid four common pitfalls that are regularly seen:

Outdated datasets Two datasets and their derivatives, DARPA-98 and KDD-99, have been extensively used to benchmark network intrusion detection models [104], with all anomaly-based techniques discussed in a recent survey [13] with reported detection rates on U2R and R2L attacks relying on either of them. However, both datasets are now more than 20 years old and have been pointed out as significantly flawed and prone to give overoptimistic results [105].

Lack of attack class distinction Most intrusion datasets include attack events from both low volume access attack classes such as R2L (Remote-to-Local) and U2R (User-to-Root) as well as attacks like DoS or port scans which generate a large number of events. Recorded events are labelled individually, creating an imbalance in the number of DoS and probing events compared other attack classes. For instance, 96% of the attack events in the CICIDS-17 dataset [84] consist of DoS and probing events. If reported detection rates do not distinguish between different attacks or attack classes, performance metrics will be dominated and potentially inflated by DoS and probing attacks.

Arbitrary false positive rates There is no agreed upon value for a suitable false positive rate in network intrusion detection. This leads many authors to report very high detection rates at the expense of having unrealistically high false positive rates, often around 5% and above. In our evaluation, we report overall AUC scores, which describe the separation of benign and anomalous traffic.

Lack of long-term evaluation To be effective, an intrusion detection system has to produce a consistently low false positive rate in the presence of concept drift. A crucial aspect when assessing the deployability of an intrusion detection system is the long-term stability of a trained model [106], which is often neglected in the literature. We include a dataset focused on long-term traffic evolution in our evaluation to demonstrate the stability and deployability of our model.

5.4 Design

5.4.1 Session construction

The raw input data, in the form of network flows, contains unordered traffic from and to all hosts in the network. To order the raw network flows, we first gather all outgoing and incoming flows for each of the hosts selected for examination according to their IP address.

The traffic a host generates is often seen as a series of *session*, which are intervals of time during which the host is engaging in the same, continued, activity [107]. In our context, flows that occur during the same session can be seen as having strong short-term dependencies. We therefore group flows going from or to the same host to sessions using an established statistical approach [107]:

If a network flow starts less than α seconds after the previous flow for that host, then it belongs to the same session; otherwise a new session is started. If a session exceeds β events, a new session is started.

We chose the number of $\alpha = 8$ seconds as we have found that on average around 90% of flows on a host start less than 8 seconds after the previous flow, a suitable threshold to create cohesive sessions according to Rubin-Delanchy et al. [107]. We introduced the β parameter in order to break up long sessions that potentially contain a small amount of malicious flows, and estimated $\beta = 25$ to be a suitable parameter. Detection rates do not seem to be very sensitive to the exact choice of β though.

A perfect session grouping would require (unavailable) information from the top layers of the network stack. We therefore use our session definition as a first approximation which we found to be useful enough for this experiment. We will discuss this issue further in Section 5.8.2 and Section 5.10.

The IAT distribution for selected hosts in the used datasets, described in Section 5.5.1, along with 90% quantile lines is depicted in Fig. 5.2.

5.4.2 Contextual modelling

Each session is now a sequence of flows that are assumed to be interdependent. We observed in an initial traffic analysis that the protocol, port, and direction of a flow as well as its size are highly dependent on the surrounding flows, which motivates their use in the modelling process. We treat flows as symbolic events that can take different states, much like words in a language model. The state of a flow is defined as the tuple

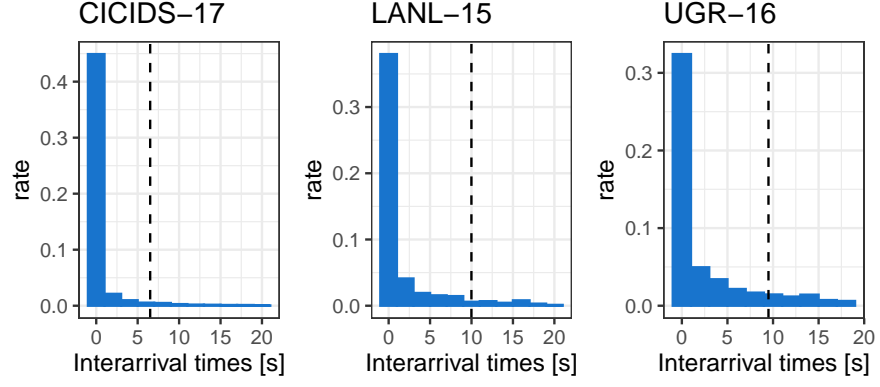


Figure 5.2: Flow interarrival distributions for selected hosts in the CICIDS-17, the LANL-15, and the UGR-16 data, with 90 percent quantile lines.

consisting of the protocol, network port, and the direction of the flow. We consider only the server port numbers, which indicate the used service, in the state-building process. We introduce the following notation:

M :	number of states
C :	number of host groups
S :	number of size groups
N_{embed}^i :	embedding dimension
N_{hidden}^i :	LSTM layers dimension
N_j :	the length of session j
$x^{i,j} \in \{1, \dots, M\}$:	the state of flow i in session j
$c^j \in \{1, \dots, C\}$:	the host group
$s^{i,j} \in \{1, \dots, S\}$:	size group of flow i in session j
$p_x^{i,j,k} = P(x^{i,j} = k j)$	the predicted probability of $x^{i,j} = k$ conditional on the other flows in session j
$p_s^{i,j,l} = P(s^{i,j} = l j)$	the predicted probability of $s^{i,j} = l$ conditional on the other flows in session j

The collection of all states is called a *vocabulary*. For prediction, the total size of a vocabulary directly correlates with the number of parameters needed to be inferred in an LSTM network, thus influencing the time and data volume needed for training.

Too large vocabularies also lead to decreased predictive performance by including rare events that are hard to predict [108]. We therefore bound the total number of states and only distinguish between the $M - 2$ tuples of protocol, port, and direction most commonly seen on a machine, with less popular combinations being grouped as “Other”. Furthermore, the end of a session is treated as an additional artificial event with its own state. The total vocabulary size is then given by M .

Our experimentation has shown that detection rates improve when including the size as an additional variable, as we discuss in Section 5.6.3. Rather than making a point estimate of the size, we want to produce a probability distribution for different size intervals. This provides better accuracy for situations in which both small and large flows have a similar occurrence likelihood. We group flows into S different size quantile intervals, with the set of all size intervals forming a third vocabulary. The $S - 1$ boundaries that separate the size intervals correspond to $S - 1$ equidistant quantiles of the size distribution in the training data.

Hosts are grouped according to their functionality (Windows, Ubuntu, servers, etc.) a distinction that can easily be performed using signals in the traffic. The group is provided to the model as an additional input parameter c^j and forms a third vocabulary.

5.4.3 Architecture selection

We now represent each session as a set of two symbolic sequences that contain between three and 27 items, in order to capture their contextual structure for the reasons described in Section 5.2. A number of techniques exist to describe such sequences, such as *Markov-models and Hidden-Markov-models*, *Finite-State-Automata*, or *N-Gram* models. However, the success of recurrent neural networks in similar applications of natural language processing over these methods suggests they would be the most appropriate architecture to capture contextual relationships between flows. In Section 5.8.3, we compare the performance of CBAM to both Markov-based models and Finite-State-Automata. Even though convolutional neural networks and feed-forward networks can be more suitable choices for specific sequential problems with tabular or regression characteristics, recurrent neural networks such as LSTMs or GRUs normally outperform them for short tokenised sequences [109]. Both LSTMs and GRUs perform similarly well and generally outperform simple RNNs.

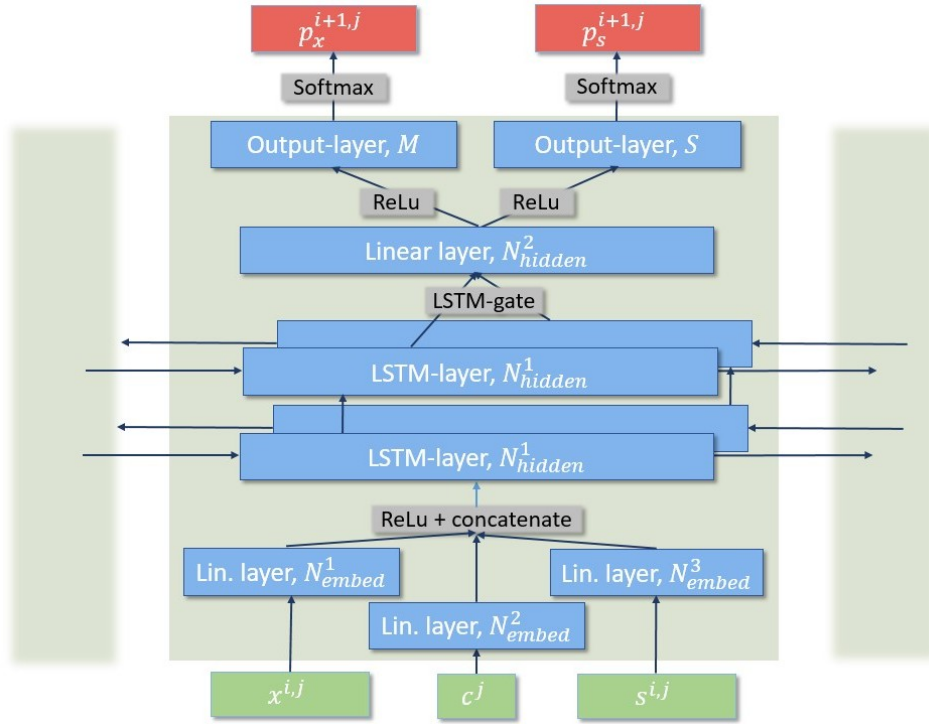


Figure 5.3: Architecture of the trained bidirectional LSTM network.

5.4.4 Trained architecture

We use a deep bidirectional LSTM network which processes a sequence in both forward and reverse direction to predict the state and size group of individual flows. The architecture of the network we trained is depicted in Fig. 5.3. The increased model complexity we present has not been explored in previous LSTM applications to network intrusion detection, and enables us to boost detection rates while lowering false positive rates, which we demonstrate in Section 5.8.

5.4.4.1 Embedding

First, each of the three vectors is fed through an embedding layer, which assigns them a vector of size N_{embed}^i , $i \in \{1, 2, 3\}$. This embedding allows the network to project the data into a space with easier temporal dynamics. This step significantly extends existing designs of LSTM models for anomaly detection and allows us to project multiple input vocabularies simultaneously without a large increase in the model size. By treating the state, the size group, and the host group as separate dictionaries, we avoid the creation of one large vocabulary of size $M \times C \times S$, which makes training faster and avoids the creation of rare states [108].

5.4.4.2 LSTM-layer

In the second step, the vectors are concatenated and fed to a stacked bidirectional LSTM layer with N_{hidden}^1 hidden cells. This layer is responsible for the transport of sequential information in both directions. The usage of bidirectional LSTM layers compared to unidirectional ones significantly improved the prediction of events at the beginning of a session and consequently boosted detection rates within short sessions, as we demonstrate in Section 5.8.1. Increasing the number of LSTM layers from one to two decreases false positive rates in longer sessions while maintaining similar detection rates, as we show in Section 5.8.2. In Section 5.10.1, we discuss why we are not further increasing the number of layers.

5.4.4.3 Output layer

The outputs from the bidirectional LSTM layers are then concatenated and fed to an additional linear hidden layer of size N_{hidden}^2 with the commonly used rectified linear activation function. We added this layer to enable the network to learn more non-linear dependencies in a sequence. We found that by adding this layer, we are able to capture complex and rare behaviours and decrease false positive rates, as demonstrated in Section 5.8.3.

Finally, the output of this layer is fed to two output layers with M and S linear output cells. These produce two numeric vectors of size M and S

$$\begin{aligned} p_x^{i,j,k}, k \in \{1, \dots, M\}, \quad & \sum_k^M p_x^{i,j,k} = 1 \\ p_s^{i,j,l}, l \in \{1, \dots, S\}, \quad & \sum_l^S p_s^{i,j,l} = 1 \end{aligned}$$

that describe the predicted probability distribution of $x^{i,j}$ and $s^{i,j}$ respectively.

The prediction loss for the state group is then given by the negative log-likelihood:

$$\text{lh}_x^{i,j} = \sum_{k=1}^M (1 - x_k^{i,j}) \cdot \log(1 - y_k^{i,j}) - x_k^{i,j} \cdot \log(y_k^{i,j})$$

with the size group loss being calculated in the same way. We calculate the total loss as the sum of the state loss and the size group loss. A visualisation of the prediction-making process is depicted in Fig. 5.4.

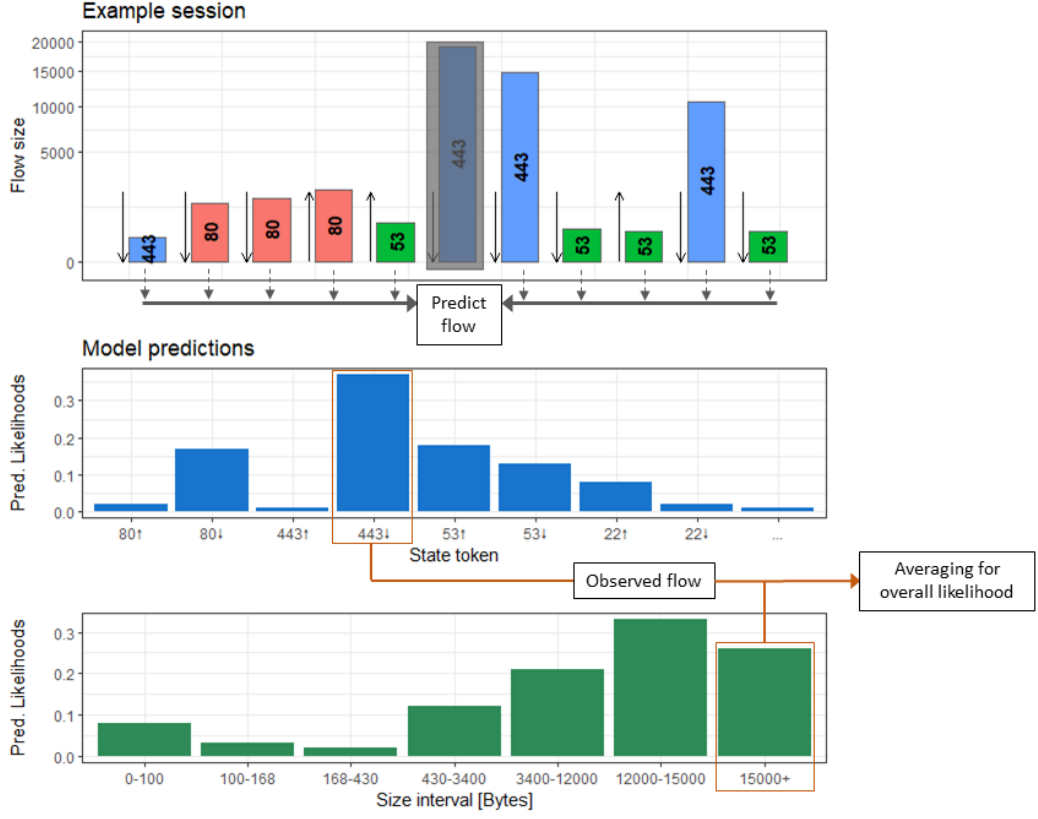


Figure 5.4: Visualisation of model prediction process.

After the training, we use the network to determine the anomaly score of a given input session via the average of the predicted likelihoods, as this measure is independent of the session length:

$$AS^j = 1 - \sum_{i=1}^{N_j} \left(\exp(\text{lh}_x^{i,j}) + \exp(\text{lh}_s^{i,j}) \right) / N_j$$

An anomaly score close to 0 corresponds to a benign session with a very high likelihood while a score close to 1 corresponds to an anomalous session with events which the network would not predict in the context of previous events. We rescale all anomaly scores for better readability, however this does not influence their ordering.

The design of the proposed anomaly score as an average likelihood has some shortcomings, which have not been considered at the time of model design and evaluation. One of these shortcomings is that while the anomaly score is designed to yield similar anomaly scores for sessions of different length, short anomalous sequences get smoothed out by other traffic the longer a session becomes as its contribution to the anomaly score is inversely proportional to the session length. Other methods such as

using the maximal value of a running average would mitigate this problem and should be investigated for future usage of CBAM.

5.4.5 Parameter selection and training

We now train CBAM and tune it to maximise its prediction performance. We train on a quad-core CPU with 3.2 GHz, 16 GB RAM, and a single NVIDIA Tesla V100 GPU, and we use minibatches of size 30 using the ADAM optimiser in PyTorch. Training a model can be achieved in under three hours.

We want to create a model that has sufficient parameters to capture complex flow dependencies, but is not overfitting the training data. For this, we split the available training data into a larger training and a smaller validation set. We then select two model configuration, one with a larger number of parameters and one with a smaller number. We then train the model for 500 epochs on the training set and observe whether the same loss decrease can be observed on the validation set. As long as the larger model is performing better than the smaller model and the validation loss is consistent with the training loss, we keep increasing the number of parameters, a standard practice to train deep learning models. The best performing parameters were $N_{\text{embed}}^1 = 10$, $N_{\text{embed}}^2 = N_{\text{embed}}^3 = 5$ for the embedding layers, and $N_{\text{hidden}}^1 = N_{\text{hidden}}^2 = 50$ for the hidden layers.

To build a more powerful model without the risk of overfitting, we use a drop-out rate of 0.5 and a weight-decay regularization of 5×10^{-4} per epoch, as suggested by Hinton et al. [110]. To increase the training performance, we use an adaptive learning of 0.0003, which decays by a factor of 2 after each fifty subsequent epochs, as well as layer normalization. The values for the learning rate and weight decay were estimated in a similar procedure as the model size.

As we mentioned above, too large vocabularies can cause problems both for model training and event prediction. We achieved the best results for $M = 200$ for the available data and computational resources. The size of the size group was chosen smaller with $S = 7$, which improves detection capabilities without increasing anomaly scores for benign sessions too much. We found that a suitable value of $C = 4$ to describe different host types, which include servers, two types of client machines depending on the operating system, and auxiliary devices (printers, IP-phones, and similar). Host groups were determined for each dataset individually and the corresponding machines labelled manually. However, for larger datasets this process is easily automated by

filtering for specific traffic events such as requests to Microsoft update servers.

	# cells	# parameters
Embedding layer	202/10/5	2055
LSTM-layer 1	50	6700
LSTM-layer 2	50	12700
Linear layer	50	2550
Softmax layer	202/10	10557
Total		34562

5.4.6 Detection method

We use a simple threshold anomaly score to identify a session as malicious. We estimate the 99.9% quantile for benign sessions in the training data, which will then act as our threshold value T . By determining T from the training data, we control the expected false positive rate in the test data. Threshold values are determined for each dataset and each host within a dataset separately.

$$T_c : P[AS_{c,j}^i \leq T_c] \leq 0.999$$

Finding an appropriate threshold value is a compromise between higher detection rates and lower false positive rates, and we chose this value to achieve false positive rates that are low enough for a realistic setting. We compare detection and false positive rates for different T in Section 5.6.1, and we give an outlook to more sophisticated detection methods in Section 5.10.

5.5 Datasets and benchmark models

5.5.1 Dataset assembly

The field of network intrusion detection has always suffered from a lack of suitable datasets for evaluation. Privacy concerns and the difficulty of posterior attack traffic identification are the reason that no dataset exists that contains realistic U2R/R2L (User-to-Root, Remote-to-local) traffic and benign traffic from a real-world environment [111]. To evaluate CBAM, we need both representative access attack traffic to

test detection rates, and background traffic from a realistic environment to test false positive rates. To ensure that both criteria are met, we selected three modern publicly available datasets that complement each other: CICIDS-17 [84], LANL-15 [38, 39], and UGR-16 [41]. The CICIDS-17 dataset contains traffic from a variety of modern attacks, while the UGR-16 dataset's length is suitable for long-term evaluation. The LANL-15 dataset contains enterprise network traffic along with several real-world access attacks.

We train models with the same hyperparameters on each dataset to demonstrate the capability of CBAM to detect various attacks and perform well in a realistic environment.

CICIDS-17: This dataset [84], released by the Canadian Institute for Cybersecurity(CIC), contains 5 days of network traffic collected from 12 computers with attacks that were conducted in a laboratory setting. The computers all have different operating systems to enable a wider range of attack scenarios. The attack data of this dataset is one of the most diverse among NID datasets and contains SQL-injections, Heartbleed attacks, brute-forcing, various download infiltrations, and cross-site scripting (XSS) attacks, on which we evaluate our detection rates.

The traffic data consists of labelled benign and attack flow events with 85 summary features which can be computed by common routers. The availability of these features makes it suitable to evaluate machine-learning techniques that were only tested on the KDD-99 data.

The benign traffic is generated on hosts using previously gathered and implemented traffic profiles to make the traffic more heterogeneous during a comparably short time span, and consequently closer to reality. For our evaluation, we selected four hosts that are subject to U2R and R2L attacks, two web servers and two personal computers.

This dataset is generated in a laboratory environment, with a higher proportion of attack traffic than is normally encountered in a realistic setting. Consequently, we need to test on traffic from real-world environments to prove that CBAM retains its detection capabilities and low false alert rates.

LANL-15 dataset: In 2015, the Los Alamos National Laboratory (LANL) released a large dataset containing internal network flows (among other data) from their corporate computer network. The netflow data was gathered over a period of 27 days with about 600 million events per day [38, 39].

In addition to large amounts of real-world benign traffic, the dataset contains a set of attack events that were conducted by an authorised red team and are supposed to re-

seem remote access attacks, using mainly the *pass-the-hash* exploit. We selected this dataset to demonstrate that CBAM is able to detect attacks in a realistic environment with low false alert rates. We isolated traffic from ten hosts, with two being subject to attack events. Two of these hosts resemble server behaviour, while the other eight show typical behaviour of personal computers.

The provided red team events are not part of the network flow data and only contain information about the time of the attack and the attacked computer. Furthermore, not all of the attack events are conducted on the network level, so it is impossible to tell exactly which flows correspond to malicious activity and which do not. Therefore we labelled all flows in a narrow time interval around each of the attack timestamps as possibly malicious. As these intervals are narrow, identified anomalies likely correspond to the conducted attack.

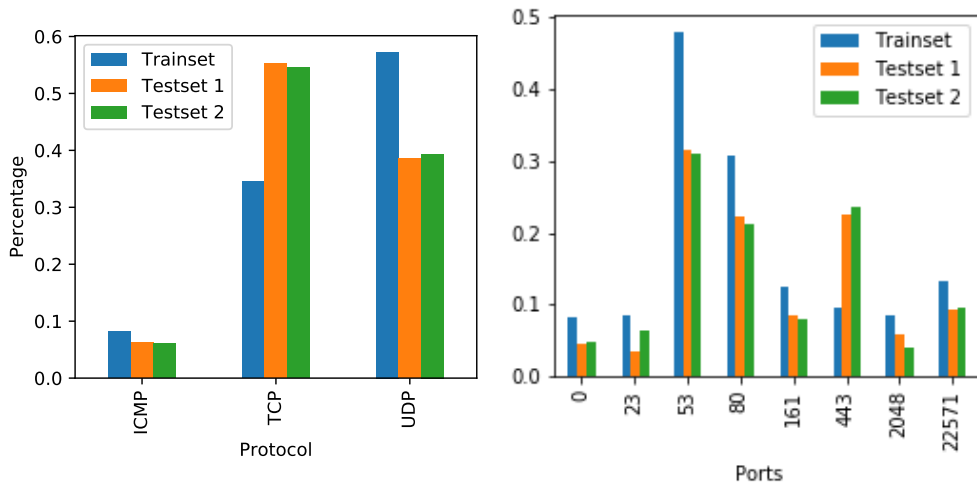


Figure 5.5: Temporal change in protocol and port usage over the different train and test intervals across selected servers in the UGR-16 dataset.

UGR-16 dataset: The UGR-16 dataset [41] was released by the University of Grenada in 2016 and contains network flows from a Spanish ISP. It contains both clients’ access to the Internet and traffic from servers hosting a number of services. The data thus contains a wide variety of real-world traffic patterns, unlike other available datasets. Additionally, a main focus in the creation of the data was the consideration of long-term traffic evolution, which allows us to make statements about the robustness of CBAM to concept drift over the 163 day span of the dataset. For our evaluation, we isolated traffic from five web-servers that provide a variety of services.

5.5.2 Dataset split

We split our data into a test set and a training set. To resemble a realistic scenario, the sessions in the training data are from a previous time interval than sessions in the test data.

To evaluate detection rates on the CICIDS-17 data, we selected the four hosts in the data that are subject to remote access attacks, two web servers and two personal computers. We choose our test set to contain the known attack data while the training data should only contain the benign data. Due to the short timespan of the dataset, we have to train on traffic from all five days, with the test data intervals being placed around the attack. In total, the test set contains 14 hours of traffic for each host while the training set contains 31 hours of traffic. While the test set for the CICIDS-17 data covers a shorter timespan, it contains more traffic due to voluminous brute-force attacks.

For the LANL data, the test set stretches approximately over the first 13 days with the training data spanning over the last 14 days. The unusual choice of placing the test set earlier the training set was made because the attacks occur early in the dataset. However as the training and test are contained in two non-overlapping intervals, a robustness evaluation is still possible.

Dataset	hosts	sessions in training set	sessions in test set	length
CICIDS-17	4	24128	32414	5 days
LANL-15	10	89480	76984	27 days
UGR-16	5	65000	480018	163 days

Table 5.2: Summary of the amount of traffic extracted from each dataset.

To test long-term stability and robustness of CBAM against concept drift, we split the UGR-16 data into one training set interval and two test set intervals, for which we can compare model performance. The training set interval stretches over the first month, with the first test set interval containing the sessions from the following two months, and the second test set interval containing the last two months. We then isolated traffic from five web-servers that provide a variety of services that show behavioural evolution. Fig. 5.5 depicts the changes of these servers in terms of protocol and port usage over the different intervals.

We chose our training data to contain about 10 000 sessions per host if possible. A summary of the amount of data in the training and test data for each dataset can be found in Table 5.2.

5.5.3 Sample imbalance and evaluation methodology

Most NID datasets include attack events from both low volume access attack classes as well as attacks like DoS or port scans which generate a large number of events. If reported detection rates do not distinguish between different attacks or attack classes, performance metrics will be dominated and potentially inflated by DoS and probing attacks. Similarly, detection results are often given in terms of precision and recall or F-measures, which are sensitive to the specific dataset balance of the majority and minority classes in a dataset. Since the ratio of attack traffic is inflated in general NID-datasets, these measures are not suitable for model comparison.

We evaluate CBAM using simple detection true positive and false positive rates, which are independent of the dataset balance. We also distinguish detection rates for different attacks, and assess overall performance by averaging these rates over attack classes rather than overall number of attack events. Since there is no agreed upon value for a suitable false positive rate in network intrusion detection, we compute ROC-curves to display the detection rates in dependence of the false-positive rate and report overall AUC-scores (*Area under Curve*), which describe the separation of benign and anomalous traffic. We use these for comparison with other models, as this measure is fairer than point comparisons. The evaluation procedure is supported by several NID evaluation surveys [112, 113].

Some researchers propose cost-based evaluation metrics by assigning false alerts and missed intrusion attempts a cost-value and tuning the detection threshold to minimise the expected cost, such as done by Ulvila and Gaffney [112]. Such a metric is however strongly dependent on the observed ratio of attack to benign traffic, which is strongly inflated in NID-datasets, and requires operational information to assign costs to a false alert or intrusion. This evaluation works well in specific cases such as DoS-attacks or cryptojacking, where server-downtime costs and attack volume are generally quantifiable, but is not applicable in cases where this information is not available or well defined [113].

Training classifiers on imbalanced dataset can affect their performance, both due to the imbalanced ratio of attack to benign traffic and the imbalance between several

attack classes. Some methods have been proposed to augment or synthetically inflate minority samples for attack traffic [114]. As an anomaly-detection method, CBAM is however trained in a self-supervised way strictly on benign traffic, with no attack traffic being present in the training data. The training stage is therefore independent of the minority class ratio in a given dataset and does not require specific balancing methods. In the evaluation stage, the above described steps apply for both classification- and anomaly-detection-based methods.

5.5.4 Benchmark comparison models

We compare our detection and false positive rates against three network anomaly models that we have re-implemented and re-evaluated.

A recent and well cited survey by Nisioti et al. [13] identified the **UNIDS model by Casas et al.** [115] as achieving the highest detection rates of remote access attacks on the KDD-99 dataset, so we chose to include this method as our first benchmark. The authors rely on subspace-projection and density-based clustering (DBSCAN) for outlier detection, and achieve detection rates of access attacks at around 80 – 85% on the KDD-99 dataset with a false positive rate of 3.5%.

Niyaz et al. [116] present a more recent deep-learning model combines anomaly detection and classification by using sparse autoencoders and detection through reconstruction error. The authors classify individual flows and claim a detection precision of 99% with a recall of 97.5%, even higher than the UNIDS model.

Finally, **Radford et al.** [96] predict sequences of individual flows between pairs of hosts using a two-layer LSTM network. Flows are tokenised according to their protocol and size, and the model detects 60% of the attacks at a false positive rate of about 2% on the CICIDS-2010 dataset. This model is closest to ours in terms of contextual anomaly detection from flow metadata, and achieves the best results out of the three benchmark models during our evaluation. We include it to highlight the improvements our design choices yield over other contextual LSTM-models.

Lastly, we include a more **shallow version of our model**, depicted in Fig. 5.6, to highlight the benefits of a deeper structure. This model only contains one LSTM-layer, and no linear layer before the output layer.

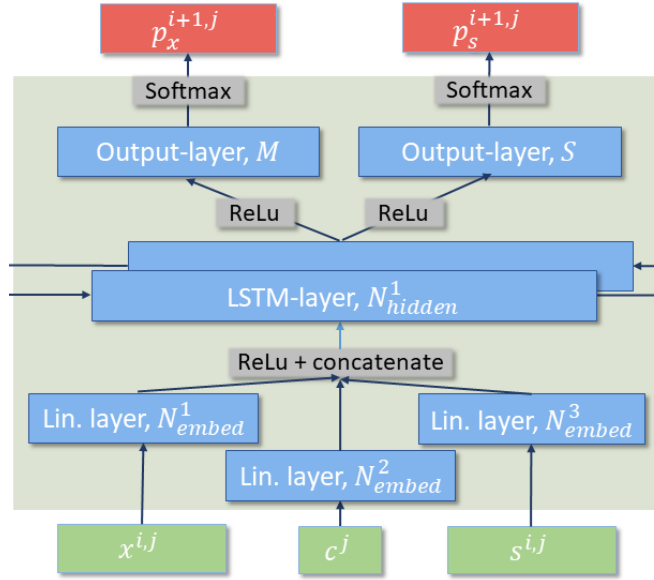


Figure 5.6: Architecture of the shallow LSTM-model version we use as a benchmark.

5.6 Detection performance

We now demonstrate that we can build an accurate and close-fitting model of normal behaviour with CBAM. We train models for each dataset separately, but without any change in the selected hyperparameters, i.e. number of hidden cells, vocabulary size, learning rate etc.

As described above, we estimate detection rates using traffic of various remote access attacks in the CICIDS-17 dataset. Table 5.3 describes the number of sessions present for each attack class.

	FTP-BF	SSH-BF	Web-BF	SQL-Inj.	XSS	Heartbleed	Infiltr.
# Sessions	243	210	88	8	41	4	17

Table 5.3: Number of sessions for each attack class in the CICIDS-17 dataset

5.6.1 CICIDS-17 results

Table 5.4 and Fig. 5.7 depict anomaly score distributions and detection rates for traffic from seven different types of attacks.

Most notable is that scores from all attacks except cross-site scripting (XSS) are significantly higher distributed than benign traffic, with median scores lying between 0.75 and 0.89. Detection rates with our chosen threshold of $T = 0.77$ are highest for

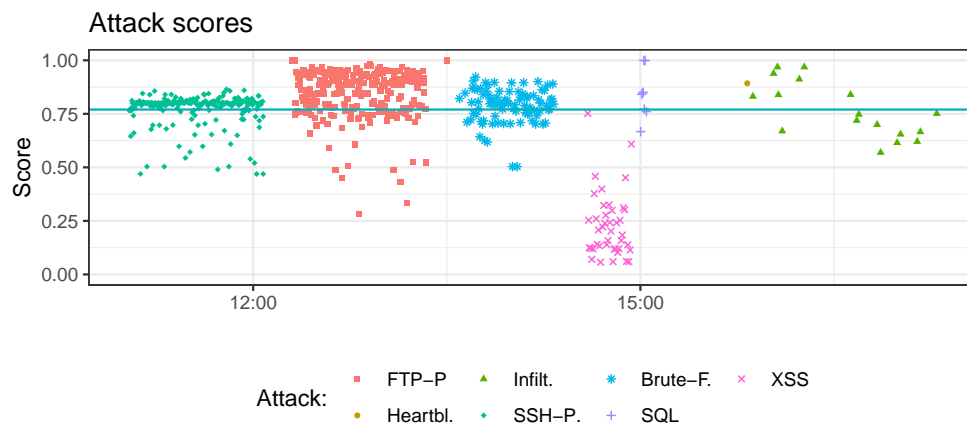


Figure 5.7: Score distribution for access attacks contained in the CICIDS-17 dataset.

Heartbleed attacks (100%), followed by FTP and SSH brute-force attacks and SQL-injections, where 91%, 74%, and 75% of all affected sessions are detected. Detection rates are lowest for XSS and Infiltration attacks. The overall detection rates we achieve are in a similar range as most unsupervised methods in Nisioti et al.’s evaluation [13], but with significantly better false positive rates.

	Anomaly scores (T=0.77)			Detection rates [%]	Shallow LSTM
	min	max	median		
Brute-force Web	0.50	0.92	0.80	0.66	0.28
FTP-Patator	0.28	1.00	0.82	0.91	0.38
Heartbleed	0.89	0.89	0.89	1.00	0.0
Infiltration	0.57	0.97	0.75	0.41	0.0
SQL-injection	0.67	1.00	0.84	0.75 2	0.21
SSH-Patator	0.47	0.86	0.80	0.74	0.67
XSS	0.06	0.75	0.20	0.00	0.0

Table 5.4: Anomaly score distributions and detection rates at threshold T for known-malicious sessions in the CICIDS-17 dataset, as well as detection rates for a less complex benchmark model (Fig. 5.6).

XSS and infiltration attacks cause the victim to execute malicious code locally. Heartbleed and SQL injections on the other hand exploit vulnerabilities in the communication protocol to exfiltrate information, and are thus more likely to exhibit unusual traffic patterns, visible as excessively long SQL-connections or completely isolated TCP-80 flows for SQL-attacks, or unusual sequences of connections initiated by the

attacked server during Heartbleed attacks.

Brute-force attacks on the other hand cause longer sequences of incoming connections to the same port of a server, in this case to port 21 for FTP, 22 for SSH, and 80 for web brute-force. Especially for port 80, such sequences are not necessarily unusual, which explains the difference in detection rates between web brute-force, which CBAM does not detect reliably, and FTP and SSH brute-force, which are detected at a higher rate. Depending on how much benign traffic the particular sessions are overlaid, the estimated anomaly scores can vary. Brute-Force attacks are not low in volume, and spread over many sessions since we introduced a maximum session length. For these types of attack, CBAM therefore only has to flag a smaller percentage of malicious sessions the attack generates to detect anomalous behaviour.

Fig. 5.8 provides *ROC* (Receiver operating characteristic) curves for each attack type. As seen, for Heartbleed, FTP brute-force, SQL injection, and infiltration attacks, CBAM starts detecting attacks with close to zero false positives.

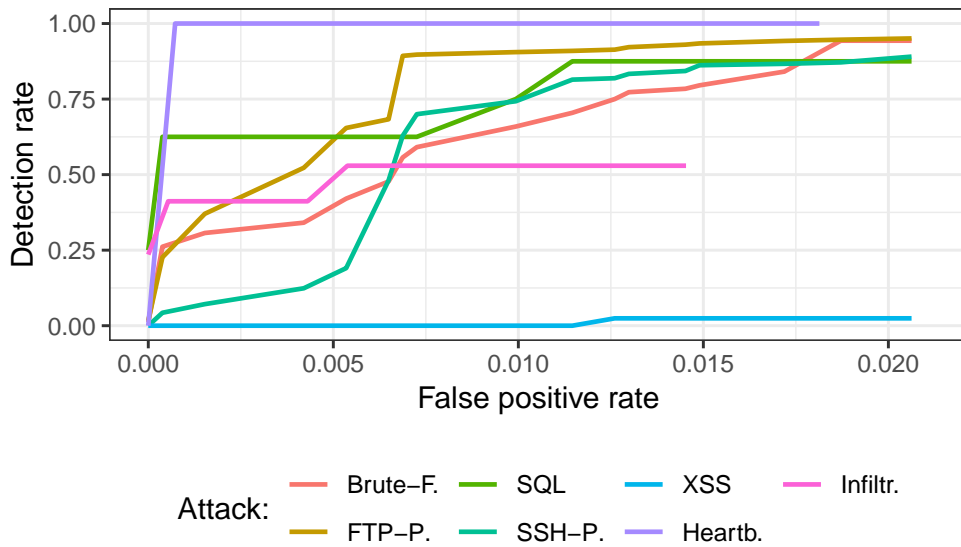


Figure 5.8: ROC-curves for different attack types in the CICIDS-17 dataset.

5.6.1.1 Comparison models

We now compare detection rates between our model and the three models described in Section 5.5.4 that we chose as benchmarks.

All three models ultimately detect anomalies when an anomaly score exceeds a threshold, which controls the balance between low false positive rates and high detection rates and usually depends on the given data. To create a fair comparison, we chose

threshold providing similar false positive rates of 0.01%, e.g. the 99.9% anomaly score quantile of the training data, which is necessary for assessing suitability for real-world deployment as we have argued in Section 5.4.6.

	1-AUC scores				
	Our model	UNIDS	Radford	Niyaz	shallow m.
Brute Force Web	0.016	0.49	0.027	0.32	0.048
FTP-Patator	0.0025	0.011	0.0048	0.16	0.0052
Heartbleed	0.0003	0.0057	0.032	0.077	0.012
Infiltration	0.046	0.033	0.35	0.15	0.11
SQL-injection	0.005	0.44	0.497	0.39	0.019
SSH-Patator	0.009	0.013	0.035	0.011	0.005
XSS	0.127	0.02	0.03	0.16	0.13
Average	0.044	0.144	0.135	0.18	0.091

Table 5.5: 1-AUC scores for our model and the implemented comparison models on the CICIDS-17 dataset. Fat numbers indicate the best value for each attack.

To assess the overall separation between benign and malicious traffic for each model, we calculated $1 - AUC$ (*Area under ROC curve*) scores by varying the thresholds for each model, and calculating the area under the ROC-curve, depicted in Fig. 5.8 and Table 5.5.

In comparison to the other benchmark models, our shallow model is capable of making some detection at the chosen false positive rate, but cannot reach the levels of our deeper model. While brute-force attacks are still detected, more nimble attacks such as Heartbleed or XSS are less distinctive from benign traffic for the shallow model. It is remarkable that by adding the described additional layers, we are able to more than double the overall detection power, as indicated by the 1-AUC-scores.

5.6.2 LANL results

We now examine whether CBAM is able to detect actual attacks in real-life traffic from the LANL-15 dataset.

As described in Section 5.5.1, we do not have labels for malicious flows in the LANL-15 data. Instead, attacks are described by narrow intervals surrounding conducted malicious activity. These intervals inevitably contain benign activity too. However, as the intervals are narrow and we saw that benign sessions only rarely receive

high anomaly scores, a session with a high anomaly score is likely to be associated with a malicious event. Of the hosts in the dataset we selected for evaluation, hosts C2519 and C754 are subject to the red team attacks. The red team activity is spread over three attack intervals **A1**, **A2**, and **A3**.

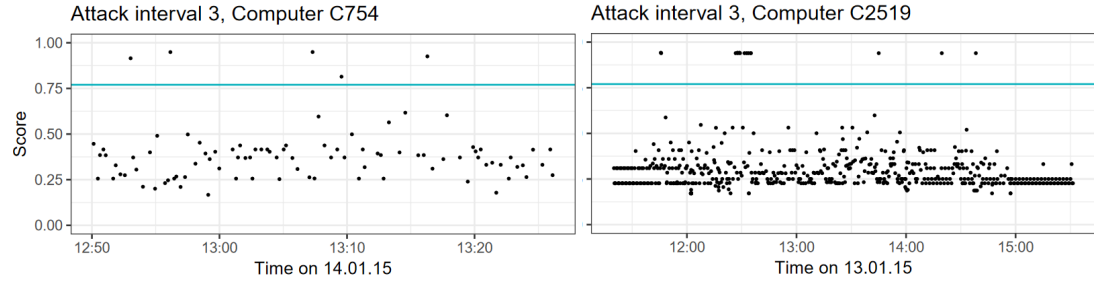


Figure 5.9: Computed scores for the third attack interval in the LANL data, along with our detection threshold.

Sessions in **A1** and **A2** have similar scores as other benign traffic, with no sessions receiving remarkably high scores. It is both possible that CBAM did not identify the malicious traffic, or that the activity in these intervals was purely host-based and did not generate any traffic.

Interval **A3** is more interesting, containing 15 sessions for host C2519 and 5 sessions for host C754 that have high anomaly scores, as depicted in Fig. 5.9.

For host C2519, each session with a high anomaly score consists of a single TCP-flow on port 445, which is usually reserved for a Microsoft SMB service. The anomaly of these sessions becomes apparent when we compare them to other sessions that contain TCP-flows on port 445.

	Proto	SrcAddr	Sport	DstAddr	Dport	Prob.End
1	tcp	C20473	N345	C2519	445	0.03
2	tcp	C2519	445	C20473	N345	0.55

Table 5.6: Exemplary session with SMB-traffic, host C2519, along with the estimated probability of the session to end.

All other sessions contain at least two subsequent flows. The model, in expectation of other following flows, assigns a very low probability to the sessions ending after a single flow. Since the analysis of the identified sessions supports their anomaly score, we believe it is very likely that these events correspond to the conducted malicious activity.

5.6.3 How attacks affect flow structures

We now examine in more detail why modelling sequences of flows is effective to detect access attacks, and how these attacks alter common flow structures. Unfortunately, the CICIDS-17 dataset, and to our knowledge all other NID-datasets, do not contain sufficient ground truth information about included attack traffic, so this analysis is based on empirical domain-knowledge of similar attacks as well as the traffic itself.

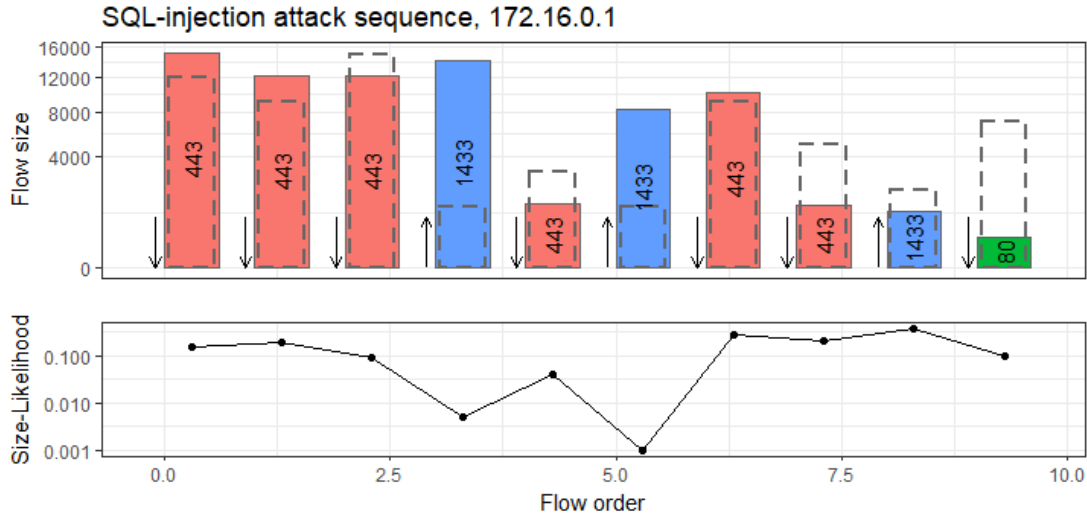


Figure 5.10: Flow-sequence in an SQL-injection attack with predicted size likelihood (log-scale). Arrows indicate flow directions (down=incoming, up=outgoing).

Fig 5.10 shows a session in the CICIDS-17 data that corresponds to a SQL-injection attack on host 172.16.0.1, a Ubuntu web server. Depicted is the order of the flows along with their direction, the destination port and the size of the flow. Dashed rectangles indicate the most likely flow size as predicted by CBAM. Below are the likelihoods of the actually observed flow sizes on a log-scale, which determine the anomaly-score of the session.

SQL-requests from a web-server typically consist of verification of user credentials or the retrieval of specific content on a web page. In an SQL-injection, SQL-code is injected into a HTTP-request that forces the server to retrieve, modify or forward additional content from an SQL-database, which can significantly increase the size of the corresponding SQL-request.

The sequence of flows in Fig. 5.10 overall resembles regular incoming HTTP requests accompanied by corresponding outgoing SQL-requests from the server. However, Fig. 5.10 clearly shows that the sizes of two of the SQL-connections on port 1433

are magnitudes larger than CBAM is predicting based on the context of the surrounding flows, which is likely caused by the injection attack. This results in a very low likelihood of the observed flow sizes and a high anomaly score for the whole session.

Fig 5.11 depicts a session that corresponds to an infiltration attack on host 192.168.10.8 in the CICIDS-17 data. Again, the figure depicts the order, direction, size and destination port of the flows, along with predictions of the most-likely sizes (dashed rectangles) and the overall likelihood of the actually observed flows.

This sequence does not resemble regular behaviour typically encountered on this host. DNS flows on port 53 are typically followed by HTTP flows on port 80 or 443, to which the model assigns a very high likelihood after the first 4 flows. However, this session contains excessively many consecutive DNS flows, which are interrupted by only one HTTP flow. Correspondingly, the likelihood for the excessive DNS flows as well as the overall session likelihood is low.

It is not completely clear how the infiltration attack triggers this abnormal behaviour. Possibly, the infiltration software is trying to retrieve the current address of a C&C-server via DNS.

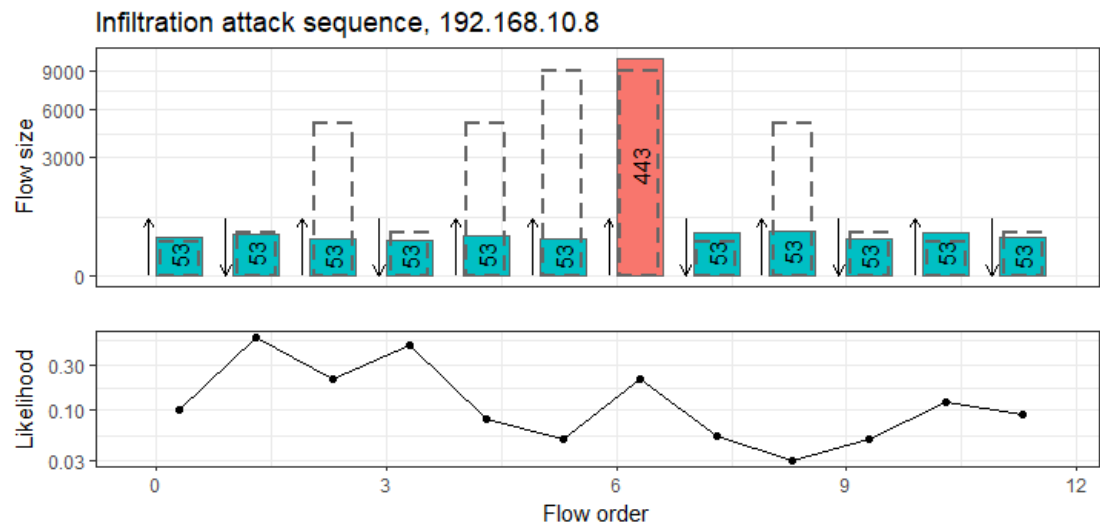


Figure 5.11: Flow-sequence in an infiltration attack with port-direction likelihood (log-scale).

5.6.4 Runtime performance

CBAM contains around 35000 parameters, which is relatively lightweight for deep learning models. The processing of a session of ten flows takes around 23ms on our

setup, which is far shorter than the average length of 5.6s of a session. In a similar comparison, our setup can process one day of activity ($\approx 15\,000$ sessions) of a web server in the UGR-16 dataset in 95s.

Considering these runtime numbers, the necessary rate of recorded flows to overwhelm our setup would need to exceed 434 flows/second. The largest rate observed for Brute-Force attacks in the CICIDS-17 dataset is 23 flows/second.

5.7 Benign traffic and longterm stability

5.7.1 UGR-16 data

We conduct the main validation of the longterm stability of CBAM on benign traffic in the UGR-16 dataset, which contains real-world traffic and spans several months. For this, we split the test data into two disjoint sets that span from May-July and August-September while being separated by one month. We then look at the quantiles and visual distribution of session scores in each test set and assess whether the score distributions and number of false positives changed as evidence on concept drift in the traffic. Fig. 5.12 depicts the score distribution of benign sessions for each dataset in the corresponding test sets.

As visible in the plot, the centre of the score distribution is concentrated very well in the lower region of the $[0, 1]$ interval, with about 50% of all sessions receiving scores in the region between 0.1 and 0.25. High scores are rare, with only very small percentages exceeding our chosen detection threshold of $T = 0.76$.

This is also reflected by the corresponding table that describes score distributions for all 5 hosts in the UGR-16 data. On average less than 0.15% of all assumed-benign sessions exceed the threshold, which would translate to fewer than ten false-alerts over the span of four months on a host with similar activity rates.

Differences in the score distributions for the two test sets are quasi non-existent. The core of the distributions are very stable, with the score quantiles differences being less than 0.03. There are some differences in the observed false positives, but the available sample size is not large enough to make any statements on any systematic differences.

A clear banding structure is visible in the plotted distributions, with most session scores being very concentrated on narrow intervals. These scores represent frequently reoccurring activities that generate very similar traffic sequences. Fig. 5.12 shows that

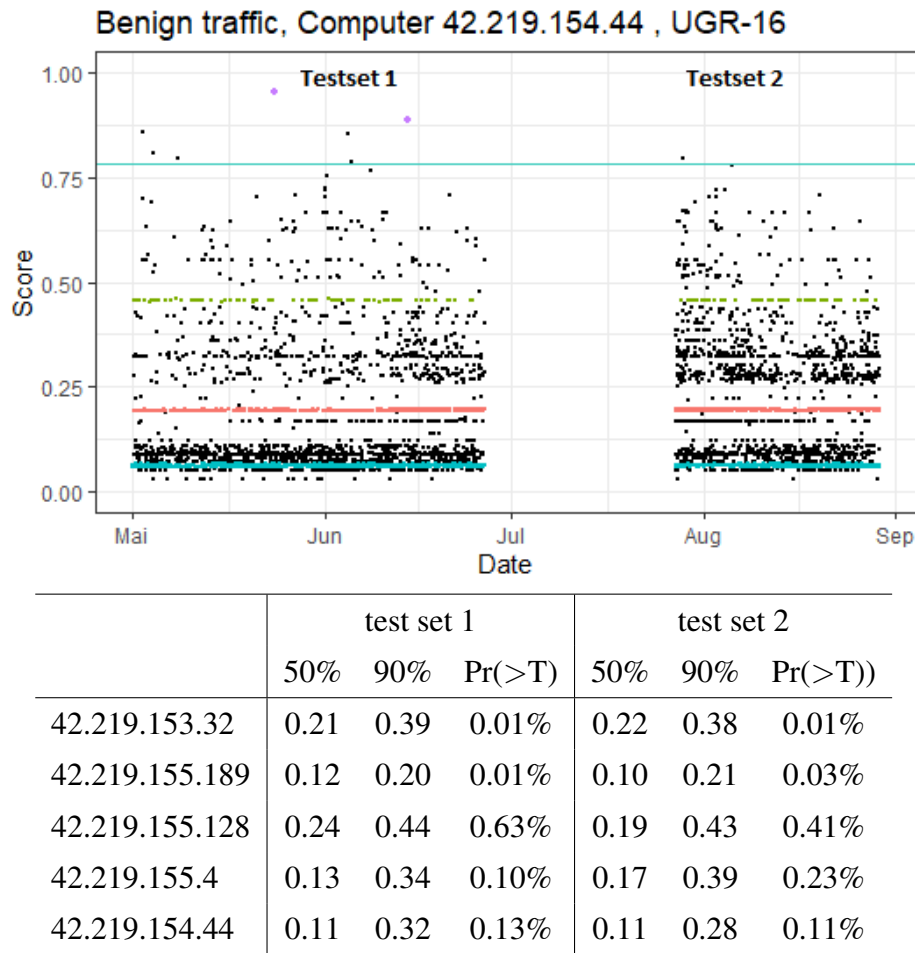


Figure 5.12: Anomaly score distributions for benign traffic in the UGR-16 data, along with an exemplary distribution plot for a selected host.

these banding structures remain virtually unchanged over several months and carry over from test set 1 to test set 2.

We coloured three of these bands in Fig. 5.12 at different levels as well as two of the observed false-positives, which we are now examining closer. Fig. 5.13 depicts the corresponding dominant session pattern that is present in each band along with the predicted likelihood for each flow. Again, the figure depicts the order, direction, size and destination port of the flows, along the overall likelihood of observed flows. For clarity, we omitted the predictions of the most-likely flow sizes (dashed rectangles).

The two lower bands, blue and red located at $AS = 0.061$ and $AS = 0.18$, represent simple and frequent HTTP- as well as corresponding NoSQL-requests and SSH activity by the server. These sessions are therefore predicted with high accuracy.

The green band at $AS = 0.45$ represents more complex and longer sessions that

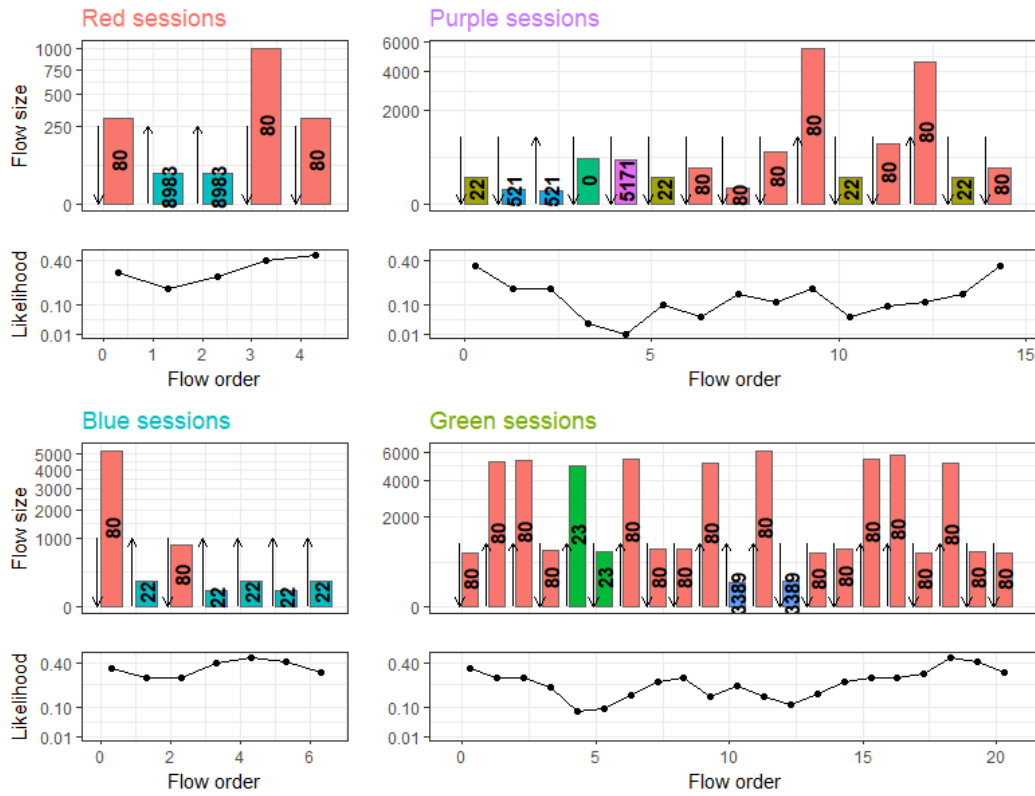


Figure 5.13: Sessions corresponding to score banding structures in Fig. 5.12, with predicted likelihoods (log-scale).

involve both incoming and outgoing HTTP-connections as well as TelNet and RDP connections. The size and order of the flows in these sessions is less deterministic than the activity in the red and blue bands. This activity is also less frequent, which explains the less accurate predictions by CBAM. The model however still recognises these sessions and is able to predict flow state and size with a non-vanishing probability, which keeps the overall session score bounded.

The two purple-coloured sessions likely represent server inspection activity, involving activity on port 0, SSH-sessions and activity on uncommon ports. This type of activity is very rare on this server and appears less deterministic than other more common activity. CBAM therefore fails to recognise the session structure and is not able to assign non-vanishing probabilities to several flows, which decreases the overall session likelihood and results in a high anomaly-score. We are not aware how often server are subject to inspections and whether this would present a problem in operational deployment. However, it seems feasible that resulting false-alerts could be linked to this administrative activity automatically or by a security analyst.

5.7.2 CICIDS-17 and LANL-15 results

We now look at the structure and stability of anomaly scores for benign traffic in the LANL-15 and CICIDS-17 datasets. The plots and tables in Fig. 5.14 depict the score distribution of presumably benign sessions in both datasets as well as describe the 50% and 90% quantiles and false-positive rates for each host. Again, score distributions for both datasets are concentrated well in the lower region of the $[0, 1]$ interval. For both datasets, the median lies between 0.06 and 0.29.

For the LANL-15 data, we observe the same banding structure as in the UGR-15 data, with most sessions being concentrated in these bands. This banding is however far less pronounced in the CICIDS-17 data, with the majority of session scores here being dispersed to a greater extent. This suggests that the traffic generation process for this dataset relies far less on reoccurring rigid activities than we observe in real-life data, which however does not seem to deteriorate the prediction performance of CBAM.

5.7.3 Importance of training data size

Host C13845 in the LANL-15 and host 192.168.10.51 in the CICIDS-17 data are an exception from the above observations, with their median anomaly score being 0.38 each and their estimated false-positive rates being 0.4% and 0.42%, which significantly exceed the average of 0.1%.

When examining host 192.168.10.51, we notice that it produced less traffic than other hosts in the CICIDS-17 data. Due to this fact, the training dataset only contains 3096 sessions or 36989 flows for this host, compared to about 10000 sessions or 115000 for host 192.168.10.25.

For hosts C13845, we see a similar picture. Because the host is less active than others in the dataset, the training data contains only 728 sessions or 2423 flows for this host, compared to 6013 sessions for the host with the next fewest training sessions.

This suggests that traffic on these hosts are not necessarily harder to predict for CBAM, but that the lack of sufficient training data prevents CBAM from learning traffic patterns for these two hosts effectively. To verify this, we examined how many sessions are necessary in the training phase to achieve similar false positives at a given anomaly threshold. We selected the hosts with the most sessions for each dataset, and reduced the number of training sessions from 10000 to 3000 and 1000. We then trained models with otherwise similar settings and compared how many additional

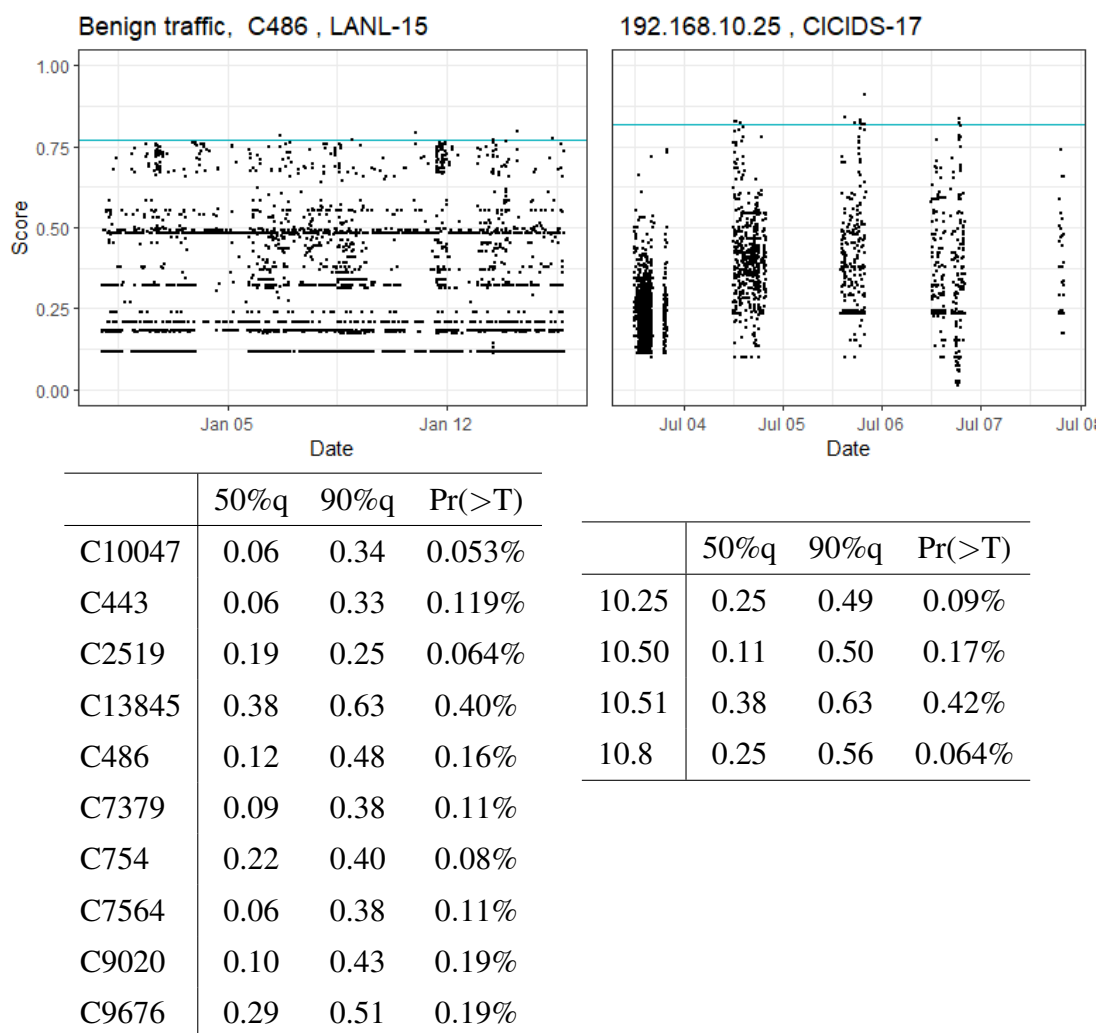


Figure 5.14: Anomaly score distributions for benign traffic in the LANL-15 and CICIDS-17 datasets.

sessions exceed the anomaly threshold. For UGR-16 and LANL-15, we examined if increasing the number of training sessions to 20000, which was not possible for the CICIDS-17 data.

Fig. 5.15 depicts the corresponding false-positive rates for each host. False-positive rates for the UGR-16 and the CICIDS-17 hosts are already significantly affected when only trained on 3000 sessions, and increase further when only 1000 host sessions are available during training. Increasing the number of sessions to 20000 however does not seem to have an effect to further improve the model.

For the host in the LANL-data, this effect is far less pronounced, and false-positives at 3000 sessions are similar to the ones at 10000 and 20000. CBAM apparently is able



Figure 5.15: Influence of number of sessions in training data for benign traffic modelling accuracy.

to learn flow predictions sufficiently from similar hosts in the dataset without depending on sessions specifically from the selected host. When we train CBAM exclusively on sessions from host C2519 without data from other hosts, the same deterioration of model prediction can be observed.

5.8 Benefit of increased model complexity

A significant part of the conducted work was concerned with improving the given network design to address insufficient predictions for several traffic phenomena and boost overall model performance. We now outline several key-steps in the design process and how they improve performance.

5.8.1 Bidirectionality for better session context

The usage of bidirectional LSTM layers compared to unidirectional ones significantly improved the prediction of events at the beginning of a session and consequently boosted detection rates within short sessions. Fig. 5.16 demonstrates this in a detailed manner: Displayed is a short session of 4 flows containing FTP and HTTP activity on host “42.219.153.32”. On the right side are the predicted likelihoods of FTP and HTTP states for each flow in the session, with the blue bars corresponding to predictions by the forward layer, while the red bars display the backwards direction and the green bars display the likelihoods after aggregating both predictions.

Fig. 5.16 demonstrates this in a detailed manner: Displayed is a short session of 4 flows containing FTP and HTTP activity on host “42.219.153.32”. On the right side

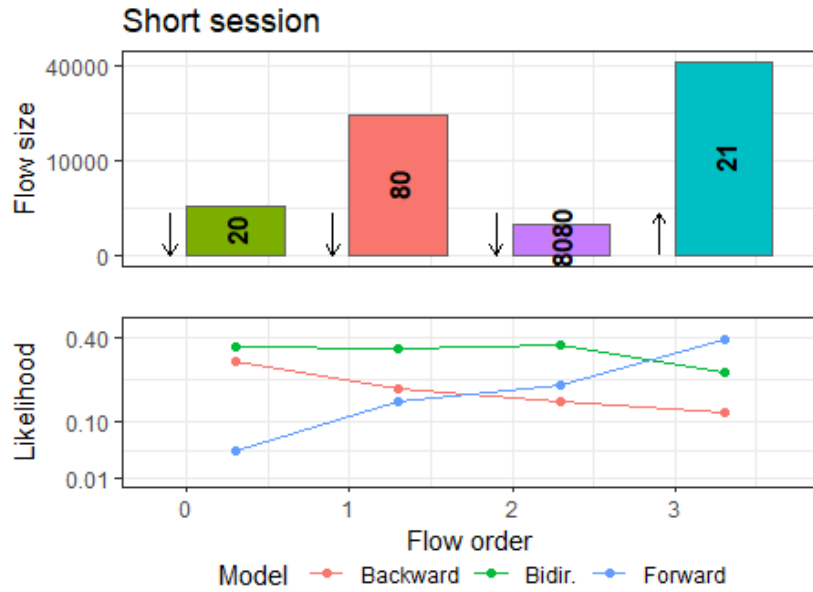


Figure 5.16: Common short session and the flow likelihoods predicted by each directional model.

are the predicted likelihoods of FTP and HTTP states for each flow in the session, with the blue bars corresponding to predictions by the forward layer, while the red bars display the backwards direction and the green bars display the likelihoods after aggregating both predictions.

When only relying on the forward direction, for the first two flows the predicted likelihoods are less than 0.07 each. The last two flows of the session are however predicted well with high likelihoods over 0.3. Because the session is short, the inaccurate predictions for the first two flows decrease overall likelihood of the session to 0.18 and the corresponding anomaly score to $AS = 0.73$, which is just below the anomaly threshold, even though this type flow sequence is quite common in the UGR-16-dataset. In a similar manner, this applies for the backward direction with the likelihoods of the last two observed flows being 0.02 and 0.03 respectively.

The cause for these phenomena is that the start of a session can differ significantly for different activities, and the LSTM-layer needs some context before recognising the specific activity and make corresponding predictions. In short sessions, the lack of accurate predictions in the first flows can then dominate the anomaly-score of the whole session.

By adding a bi-directional layer, we are able to provide context for these initial flows in a session as well by looking at later flows first. The green bars in Fig. 5.16 displays this: By basing predictions both on the output of the forward- and the backward-

		Likelihood of flows 1&2		FP rate [%]	
		unidir	bidir	unidir	bidir
UGR-16	All sess.	0.13	0.27	0.31	0.12
	sess.<5 flows	0.19	0.41	1.6	0.09
CICIDS-17	All sess.	0.09	0.29	0.37	0.18
	sess.<5 flows	0.05	0.30	1.7	0.13

Table 5.7: Average likelihood of first two flows in a session and false-positives for uni- and bidirectional model.

layer, the bidirectional model is able to predict flow likelihoods significantly better and thus assign the session a much lower anomaly-score.

Table 5.7 displays how much we could decrease false-positive rates by replacing the unidirectional LSTM layer with a bidirectional one. Overall, the false-positives decreased by 61% for the CICIDS-17 dataset, and by 52% for the UGR-16 dataset. More strikingly, when only looking at short sessions that contain less than 5 flows, we were able to reduce false-positives by 94% and 87% respectively.

5.8.2 Additional layers for complex session modelling

The inclusion of a second LSTM-layer as well a subsequent linear layer allows CBAM to capture more complex behaviour in long sessions as well as remember rare behaviour more quickly. It also increased the average predicted likelihood for flows overall.

To examine the benefit of the described model depth, we compare it to a more shallow version that lacks the second LSTM- and linear layer, as depicted in Fig. 5.6, which was trained under otherwise similar conditions. Overall detection rates for this model can be found in Table 5.4 and 5.5, while score distributions can be found further down in Table 5.8. Here, we examine in more detail how the increased model depth allows better predictions for complex sessions.

Fig. 5.17 displays two different types of activities, A and B, which are common in the CICIDS-17 data. The structure in these sessions can be observed frequently with only minor variations. Consequently, the sessions are predicted well by both the original and the more shallow model.

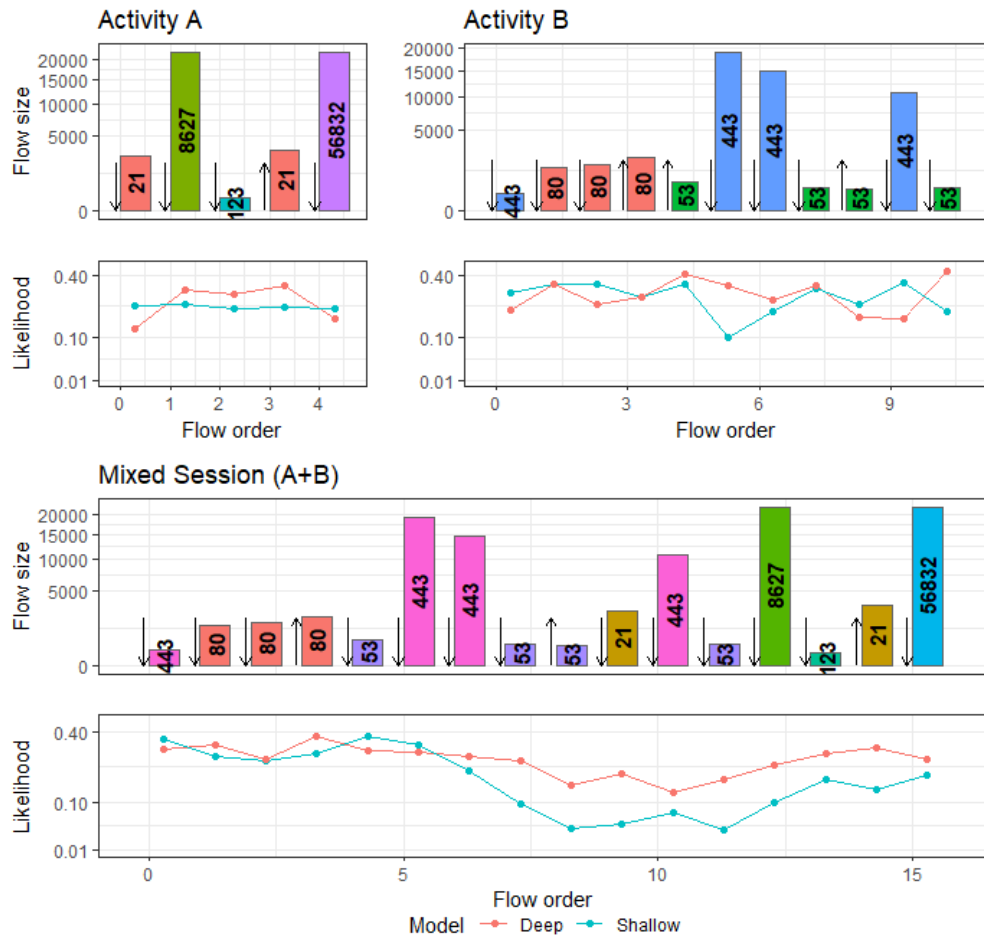


Figure 5.17: Predictions for two activities in isolated sessions and in a mixed session.

However, traffic from two or more activities can sometimes occur simultaneously and thus get grouped into the same session. Fig. 5.17 shows how the traffic from activity A and B are overlapping in a session, which makes the structure in the session more complex to predict.

The displayed likelihoods show predictions by the shallow model are accurate for flows at the beginning of the session well, but deteriorate once encountering flows from activity B. Prediction accuracy by the more complex model is also decreasing, but remaining on a sufficient level to assign this session an anomaly score of $AS = 0.51$, compared to $AS_S = 0.79$ for the shallow model. When looking at the activation in the LSTM-memory cell, we see that similar neurons as in activity A are activated at the beginning of the session, which shifts during the course of the session and resembles a more similar activation as in activity B at the end.

The improvements achieved by adding these additional layers could suggest that increasing the number of layers even further will decrease false positive rates even

further, which we discuss in Section 5.10.

5.8.3 Comparison with simpler models

In this section we aim at studying whether the higher complexity of an LSTM network is necessary for the task of detecting contextual network anomalies, or whether simpler baseline methods can achieve the same results. For comparison purposes, we implemented a first-order Markov Chain (MC) and a Non-Deterministic Finite Automata (NFA) model. Both methods are widely used in sequence modelling, and have been applied successfully to security problems [117, 118]. In contrast to LSTMs, Markov Chains have no memory past the last event while NFAs can distinguish between different types of sequences via state-merging, and give corresponding transition probabilities.

Similar to our LSTM model, Markov Chains and Finite Automata predict state transition probabilities, which is why we can employ the same anomaly score computation. However, computational costs increase quadratically with the number of states, and a separation of state vocabularies is not possible. We restrict both comparison models to the above described port-direction states. Models and detection rates were determined on the CICIDS-17 dataset. Table 5.8 shows distribution characteristics of benign and malicious sessions for our shallow LSTM-model, the Markov Chain model, and the NFA. It shows that CBAM outperforms these baseline methods, but also that the automata performs better than Markov Chains. While the Markov Chain is practically not able to make any distinction between malicious and benign traffic, the automata model shows some, albeit limited ability to identify anomalous sessions, mainly for the three types of brute-force attacks. This order shows the importance of sequence memory for contextual anomaly detection, and confirms our previous comparison of the suitability of Markov Chains and NFAs for network intrusion detection [119].

5.9 Related work

The application of recurrent neural networks to network intrusion detection has risen in popularity recently. LSTM-models for web attack detection, such as by Yu et al. [120], improve detection rates of simpler preceding models such as Song et al. [121]. They rely on deep packet inspection, and are often targeted at protecting selected web-servers rather than network-wide, due to a lack of computational scalability and in-

	shallow LSTM	Markov MC	NDFFA
Ben. 50%q	0.22	0.61	0.55
Ben. 90%q	0.55	0.81	0.86
Ben. 99%q	0.73	0.89	0.96
Mal. 50%q	0.70	0.60	0.68
Mal. 90%q	0.85	0.83	0.89
Mean AUC	0.86	0.53	0.64

Table 5.8: Score distributions for simpler models.

creasing traffic encryption. Methodologically, vocabularies are created from string sequences with well-known NLP methods, while CBAM provides a new vocabulary-construction method suitable for traffic metadata.

The majority of LSTM-based metadata approaches rely on labelled attack data for classification, and do not have the scope of anomaly-based models to detect previously unseen attacks. A prominent example of this comes from Kim et al. [103], who classify flow sequences based on 41 numeric input features. Anomaly-based approaches, such as ours, mostly rely on iterative one-step ahead forecasts, with the forecasting error acting as the anomaly indicator. This is for instance done in GAMPAL by Wakui et al. [122], who use flow data aggregation as numerical input features, which are computationally easier to process, but cannot encapsulate high-level information such as the used protocol, port, or direction. These models are best used for detecting high-volume attacks. Apart from our work, only Radford et al. [96] create event vocabularies from flow protocols and sizes. We use a more sophisticated model in terms of stacked recurrent layers and embeddings for more input features, which results in higher detection rates, as demonstrated in see Section 5.8. The HCRNNIDS model by Kahn provides an interesting adaption of hybrid convolutional recurrent networks typically used in video modelling to intrusion detection [123] with promising results. In comparison to CBAM, this model is applied to individual flow features rather than flow sequences, and is trained as a classifier rather than an anomaly-detection model.

Encoding methods are increasingly used in combination with LSTM networks to create embeddings of packet or flow sequences, such as done by Zhong et al. [124] for

anomaly detection. Zhou et al. [125] use embeddings to facilitate anomaly-detection that is robust against dataset imbalances. Liu et al. [114] use embeddings to augment and inflate minority class data samples for the same purpose.

Berman et al. [126] have surveyed recent deep-learning techniques for network intrusion detection as well as other cyber-security applications. They assess that many recurrent methods are state-of-the-art, but do not reach a conclusion whether they perform better than convolutional or generative methods.

Notable work outside of network traffic includes Tiresias [127], an LSTM model for security event forecasting with great accuracy, and DeepLog [128], an LSTM network to learn a system's log patterns (e.g., log key patterns and parameter values) from normal execution. The design of Tiresias has similarities to ours, but the scope of the model is attack forecasting rather than intrusion detection, and relies on both different input data in the form of IDS logs as well as different evaluation metrics. DeepLog is combined with a novel log parser to create a sequence of symbolic log keys, which is then also modelled using one-step forecasting. The authors achieve good detection results in regulated environments such as Hadoop with limited variety of events (e.g., 29 events in Hadoop). Here, CBAM goes further by being applied to a much more heterogeneous data source and creating a more than 30 times larger vocabulary. Han et al. [129] have recently proposed a deep graph-net based anomaly-detection method for provenance based data that demonstrates how effective neural anomaly-detection methods at detecting unknown intrusions.

5.10 Limitations and evasion

5.10.1 Limitations

CBAM is an initial application of short-term contextual modelling on network traffic that demonstrates the potential of contextual traffic models for intrusion detection. Although we use a relatively simple model with few, but carefully selected input features, we outperform sophisticated methods while retaining low false positive rates. The detection rates are to be taken with care as the available access attack data is small, synthetic and contains just a limited number of attack classes. The detection rates in the cross-evaluation on a real-world access attack in the LANL-15 data gives us confidence that CBAM's performance is reproducible in real-world scenarios, but additional data is required for an ultimate conclusion.

A frequently asked question concerns whether low false-positive rates carry over from the synthetic background traffic in datasets such as CICIDS-17 to real-world scenarios [14]. We believe that this was sufficiently demonstrated by the long-term evaluation and the observed score stability on the UGR-16 real-world dataset.

The improvements achieved by adding additional layers could suggest that increasing the number of layers even further will decrease false positive rates even further, and is certainly worth exploring in future work. However, as discussed in Section 5.7.1, the current main source for false positives are rare activity events which are not contained in the training data and are therefore not be recognised by the model. To make significant reduction in the false positive rate, we would need to train on datasets spanning more computers or over longer time periods. We are however aware of the difficulties involved in creating datasets for NIDS evaluation.

5.10.2 Evasion and resilience

Evasion tactics and corresponding model resilience against them have been a concern in the development of NIDS. We specifically focused on short-term sequential anomalies as they are often an unavoidable by-product of attack sequences, and it is thus very difficult for an attacker to perturb attack sequences that rely on a specific exploit without pre-existing control over the victim device or other network devices. We therefore believe that CBAM is relatively robust against evasion. However we identified potential improvements for future work.

A specific evasion tactic that has been discussed extensively in the context of machine learning is model poisoning in the training/retraining phase. A great difficulty for the attacker is the fact that CBAM uses sequences of symbolic events rather than continuous parameters. The introduction of a gradual shift is therefore not possible in a direct way as the alteration of individual events would look anomalous straight away. Furthermore, it is normally not possible for an attacker to alter individual events significantly without pre-existing control over network devices or specific exploits, i.e. the change of the port or size would normally cause an error in the communication. It is thinkable that an attacker could increase the predicted probability of specific events patterns more gradually by overlaying traffic stemming from 3rd party devices. However, the attacker would either need control of these devices or the ability to monitor traffic to the victim device in real-time, both of which is usually not available. We also showed in Section 5.7 that short-term contextual traffic patterns remain stable over

several months, which means that retraining of CBAM is only necessary at a low rate and attackers will have to wait for a long time to execute successful model poisoning.

An issue we encountered is the overlay of malicious and benign traffic. Currently, the existence of known traffic patterns in a session can deplete the overall anomaly score of a session. A potential evasion tactic could therefore try to conceal an attack behind benign communication on the victim device, an already common approach for C&C communication. Possible improvements for this issue are a refined notion of a session that groups related traffic better, and a better scoring method that identifies smaller anomalous sequences in an otherwise normal sequence of flows. Additionally, developing more sophisticated detection methods from the computed scores may boost detection rates.

A potential tactic that an attacker could use to make a multi-step attack consisting of multiple attack connections appear less anomalous is to artificially slow down connections. If a connection in an attack sequence is slowed down enough, a response from the victim or a follow-up connection from the attacker triggered by this connection could fall outside of the session window of 8 seconds. This way, an attacker could place each connection in an attack sequence in a separate session, where the anomaly of the connections would appear less anomalous than if most or all connections fall in the same session. Future work should therefore examine if consecutive sessions could be compared and if necessary concatenated or correlated.

5.10.3 Future Work

CBAM is an initial application of short-term contextual modelling to network traffic that demonstrates the potential of contextual traffic models for intrusion detection. Although we use a relatively simple model with few, but carefully selected input features, we outperform sophisticated methods while retaining low false positive rates. Future work may further improve detection and false positive rates to consistently detect more types of contextual anomalies. Future work to extend the input or output feature space, or to grow the size of the model will require careful design and calibration, but may potentially achieve even better results and more accurate traffic representations.

An obvious step for improvements is to extend the existing model onto more input features, making it harder for attacks to resemble normal short-term interactions in network traffic while still performing their intended activity. Features such as packet distribution features for individual flows are hard to spoof and can give more context

to common sequences.

An issue we encountered is the overlay of malicious and benign traffic. Currently, the existence of known traffic patterns in a session can deplete the overall anomaly score of the whole session. Possible improvements for this issue are a refined notion of a session that groups related traffic better, and a better scoring method that identifies smaller anomalous sequences in an otherwise normal sequence of flows. Additionally, developing more sophisticated detection methods from the computed scores may boost detection rates.

To improve predicted probabilities at the start of a session, it may be helpful to process the event sequence both from the start and the end using bidirectional modelling. False positive rates may be decreased by computing anomaly scores differently for shorter sessions.

Finally, the performance evaluation of our model would benefit from longer datasets with more comprehensive real-world attacks to assess its the deployment potential. We are however aware of the difficulties involved in creating datasets for NIDS evaluation.

All these steps will require careful design and calibration, but may potentially achieve even better results and more accurate traffic representations.

5.11 Conclusion

CBAM presents a new and promising angle to anomaly-based intrusion detection that significantly improves detection rates on the types of network attacks with the lowest detection rates. By assigning contextual probabilities in flow microstructures, CBAM improves detection rates of low-volume remote access attacks and outperforms current state-of-the-art anomaly-based models in the detection of several attacks while retaining significantly lower false positive rates. Furthermore, CBAM retains low false positive rates for periods stretching several months. Our results provide good evidence that designing anomaly detection models based on observed traffic microstructures may in the future help decrease the threat of previously unseen vulnerabilities and malware aimed at acquiring unauthorised access on a host. We specifically focused on short-term flow anomalies as they are often an unavoidable by-product of an attack thus very difficult for an attacker to avoid without pre-existing control over the victim device or other network devices.

Chapter 6

Stepping-stone detection and evasive microstructure control

6.1 Introduction

The problem of stepping-stones detection (SSD) has been studied for over 20 years, yet the body of literature fails at providing an informative overview of the detection capabilities of current methods. In this chapter, we set out to do just that by evaluating and comparing a number of selected state-of-the-art approaches on a new and independently generated dataset.

In a stepping-stone attack, malicious commands are relayed via a chain of compromised hosts, called stepping-stones, in order to access restricted resources and reduce the chance of being traced back. Real-world attacks using stepping-stone chains include Operation Aurora [130], Operation Night Dragon [131], the Black Energy [132] attack on the Ukrainian powergrid, and the MEDJACK [133] attack where medical devices were used as stepping-stones. The European Union Agency for Cybersecurity currently classifies stepping-stone attacks as one of the top ten threats to IoT-devices [134].

The detection of interactive stepping-stones is challenging due to various reasons. Packet-based methods are computationally expensive while false-positives can render a method unusable. Attackers are not constrained to specific proxy techniques and can encrypt their communication to make detection through content comparison impossible. Many methods therefore aim to detect traffic microstructures that reoccur in the relayed traffic. However, attackers can obfuscate these microstructures with evasive perturbations. Like many intrusion attacks, stepping-stones are rare and there exist no

public datasets, leading researchers to evaluate their methods on self-provided private data, which makes a direct comparison of the achieved results impossible.

This chapter discusses the following results:

1. We describe a framework to generate data that represents realistic stepping-stone data without bias to particular detection mechanisms. Our framework is scalable and capable of generating sufficient variety in terms of network settings and conducted activity.
2. We re-implemented eight SSD methods that represent the current state-of-the-art and provide a fair evaluation of their capabilities in a number of settings.
3. Our evaluation shows that while most methods can accurately detect stepping-stones through traffic microstructure propagation, detection rates plummet when appropriate perturbations are inserted. This result disproves the claims made for multiple methods that their detection rates are robust against chaff perturbations.

Thesis context: This chapter builds up on the previous chapter by displaying several methods that leverage microstructures for the detection of stepping stones, and displays how an attacker can evade detection by altering specific microstructural aspects of the generated traffic.

This chapter is mostly consisting of work published in “Evading stepping-stone detection with enough chaff” (Henry Clausen, Michael S. Gibson, and David Aspinall, 2020 [6]).

6.1.1 Outline

The rest of the chapter is organised as following: Section 6.1 provides an introduction and background to the problem of stepping-stone detection. Section 6.3 discusses the particular design of the data generation framework. Section 6.4 presents the dataset arrangement in terms of background and attack data and discusses evaluation methods. Section 6.5 discusses the selection process, properties, and implementation of the eight SSD methods that we implemented for evaluation. Section 6.6 discusses the results achieved by the implemented methods on the given data. Section 6.2 discusses related work.

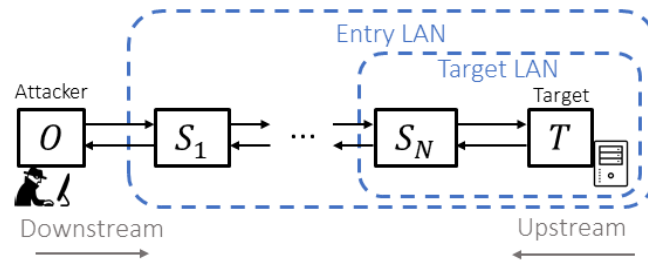


Figure 6.1: Depiction of an exemplary stepping-stone chain.

6.1.2 Background and threat model

Stepping-stones were first conceptualised by Staniford-Chen and Heberlein in 1995 [135]. In an interactive stepping-stone attack, an attacker located at the origin host, called *host O*, sends commands to and awaits their response from a target, *host T*. The commands and responses are proxied via a chain of one or more intermediary stepping-stone hosts, called *host S₁, ..., S_N*, such as depicted in Fig. 6.1. Once a host *S_i* is brought under control, it can be turned into a stepping-stone with simple tools and steps. Some of the most common set-ups are port forwarding via SSH-tunnels, setting up a backpipe with NetCat, or using Metasploit to set up a SOCKS proxy [136].

Stepping-stone detection (SSD) is a process of observing all incoming and outgoing connections on a particular host *h_i* and determining whether it is used to relay commands. This is generally done with no prior information about any other stepping-stone hosts *S₁, ..., S_N* or the endpoints *O* and *T*. A popular approach to SSD is to compare connections pairwise to identify whether they carry the same information. To avoid detection, several evasive flow transformation techniques exist that aim at decreasing observable correlation between two connections in a chain.

- **Packet transfer delays/drops:** An attacker can choose to apply artificial delays to forwarded packets, or drop certain packets to cause retransmission, in order to create temporal disparity between connections. Researchers often assume the existence of a maximum tolerable delay [137].
- **Chaff perturbations:** Chaff packets do not contain meaningful content and are added to individual connections in a chain without being forwarded. Adding chaff perturbations can be used to shape the connection profile towards other traffic types.
- **Repacketisation:** Repacketisation is the practice of combining closely adjacent packets into a larger packet, splitting a packet into multiple smaller packets, or altering the packet content to change observed packet sizes and numbers.

In our evaluation, we set out to understand the effect of different evasive methods on detection rates. We are focusing on entirely on interactive stepping-stones as described above, in which an attacker relays commands in real-time to their target. So called *store-and-forward* stepping stones, which are described in 6.2 are not discussed as their detection would follow a very different approach and requires different types of data.

Similarly, we are looking at methods that detect stepping-stones by identifying relayed interactive connections from a network-perspective. We are not looking at methods that use host log data to reconstruct connection paths through a network or that describe unusual activity patterns related to stepping-stones, as simulating temporal stepping-stone activity realistically is difficult, with little empirical data to support specific assumptions.

6.2 Related work

6.2.1 Testbeds and data

In 2006, Xin et al. [138] developed a standard test bed for stepping-stone detection, called *SST* that generates interactive SSH and TelNet connection chains with variable host numbers. In contrast to our work, the authors give little detail on implemented evasive tactics, and is not available any more.

An approach to use publicly available data comes from Houmansadr et al. [16], who simulate stepping stones by adding packet delays and drops retroactively to connections in the CAIDA data [43]. While this procedure seems sufficient for the evaluation of watermarking methods, it falls short on simulating the effects of an actual connection chain and leaves out chaff perturbations.

We find that when authors evaluate methods on self-generated data, tested evasive behaviours are often lacking analytical discussion and their implementations are too simplistic, leading to increased detection rates. An example of this can be seen in the evaluation of Ano1 [139], where a standard option in Netcat is used to generate chaff perturbations for evaluation, or for PContext [140] where simulated chaff is added randomly after the traffic collection. Furthermore, often a relatively low limit on the amount of inserted chaff perturbations is assumed without obvious reason, thus avoiding evaluation at higher ratios.

6.2.2 Detection methods and threat models

Coskun et al. [141] identify another form of stepping-stones called *store-and-forward*, which transfer data within files in a non-interactive manner. Though harder to detect than interactive connections, this procedure limits the attackers ability to explore the target, which is why SSD research has been primarily concerned with interactive stepping-stones.

A different direction in SSD that we did not discuss in this work focuses more on the general communication behaviour of selected hosts rather than individual connections. Features include the timely correlation of connecting IP-address on a selected host or unusual paths of simultaneously existing connections within a computer network. Graph-based and behavioural models such as *Hopper* by Ho et al. [142] or Apruzzese et al. [143] do not examine individual connections but instead focus on the overall temporal activity of hosts in a network or as appearing in host logs. Since our dataset focuses only on individual connections and corresponding evasion, we are not able to include graph-based or behavioural approaches in our evaluation. For the same reason, we did not include *store-and-forward* stepping-stones, which transfer data in a non-interactive manner and should thus be identified through their temporal behaviour.

6.3 Data generation setting

6.3.1 Containerisation with DetGen

To ensure reproducibility, we rely on containerisation. A container is a standard unit of software that runs standalone in an isolated user space in order to remove platform dependencies and ensure repeatability. The use of containerisation for this project follows the same traffic generation paradigm as used for our DetGen framework, which we described in Chapter 3.

6.3.2 Simulating stepping stones with SSH-tunnels and Docker

We want to capture data not only from one interaction in a fixed stepping-stone chain, but from many interactions and chains with different settings. For that, we run multiple simulations, with each simulation establishing a stepping-stone chain and controlling the interactions between host O and host T .

A simulation begins with the start-up of the necessary containers and ends with

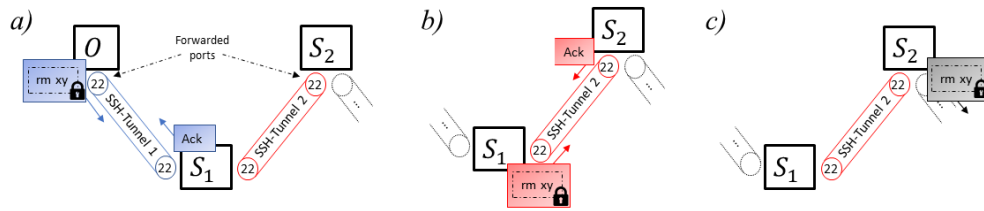


Figure 6.2: Depiction of the way a command is packetised, encrypted, and travels through the different stages of the stepping-stone chain via SSH-tunnels.

their takedown. We simulate host O , host T , and host S_1, \dots, S_n with SSH-daemon containers. To establish a connection chain, we connect these containers via SSH-tunnels, with the first tunnel forwarding a port from host O to host S_1 , which is then forwarded to host S_2 by the second tunnel etc. As mentioned by Gordon Fraser [136], this is one of the most common pivoting methods for attackers. Traffic is captured both at host T and host S_n , which acts as the final stepping-stone in the chain. Fig. 6.2 depicts a packet transfer via an exemplary chain.

6.3.2.1 Simulating interactive SSH-traffic

In order to generate enough data instances representing interactive stepping stone behaviour, we automatised the communication between host O and host T . For each simulation, we generate a script which passes SSH-commands from host O to host T .

To mimic a user's actions, we compiled a command database which consists of common commands and their usage frequency, similar to [138]. Commands are drawn randomly according to their usage frequency and concatenated to a script. Commands can either be atomic, such as "ls-la" or "pwd", or compound commands such as inputting text to a file. Command inputs are randomised appropriately when a compound command is drawn. A script ends once the *End*-command is drawn at random from the command catalogue.

To simulate human behaviour that is reacting to the response from host T , all commands are separated by *sleep*-commands for time t , which is drawn from a truncated Pareto-distribution. Paxson et al. [75] have shown that inter-packet spacings corresponding to typing and "think time" pauses are well described by Pareto distributions with a shape parameter $\alpha \approx 1.0$.

6.3.2.2 Simulating different network settings

Hosts in a stepping-stone chains can be separated by varying distances. Some may sit in the same LAN, while others may communicate via the Internet from distant geographical locations, which influences the round-trip-time, bandwidth, and network

reliability.

To retard the quality of the Docker network to realistic levels, we rely on the emulation tool NetEm, which allows users to artificially simulate network conditions such as high latency, low bandwidth, or packet corruption/drop [64]. We set the network settings and bandwidth limit for each host container individually before each simulation to allow hosts to experience different settings.

6.3.3 Evasive tactics

6.3.3.1 Adding transfer delays

To simulate evasive behaviour, we add transfer delays to forwarded packets. This method, often called *jittering*, can destroy time-based watermarks in packet flows and help decrease observable correlation between two connections. The delays are added using NetEm. We draw delays from a uniform distribution, covering the interval $[0, \delta_D]$. This particular choice has been suggested by Padhye et al. [144] in order to mimic the interarrival distributions of streaming services. The value of δ_D is fixed before each simulation and can be varied to allow for different degrees of packet jittering. We explore values for δ_D up to 1500 ms, with values above leading to unstable communication. Results in Section 6.6 show that this is enough to render watermarking methods and most flow correlation methods obsolete.

6.3.3.2 Adding chaff perturbation

We insert chaff packets without actual information to individual connections in the chain using a Netcat client. To add and filter packets in a connection, we open additional ports in each SSH-tunnel that are however not forwarded through the entire chain. Padhye et al. [144] suggest to generate chaff that mimics the flow characteristics of streaming services to both spread the added perturbations evenly across the connection and increase the difficulty of detecting the perturbation itself. For this, packet sizes are drawn from a truncated Lognormal-distribution with mean μ_C , while transmission intervals are drawn from a uniform distribution that covers the interval $[\delta_C/2, \delta_C]$ to mimic a constant packet flow. By adjusting δ_C , we can control the amount of chaff sent.

6.3.3.3 Repacketisation

By design, SSH-tunnels perform repacketisation along with re-encryption and independent packet confirmations.

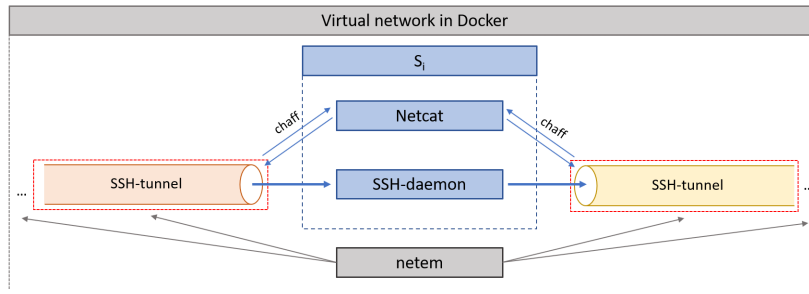


Figure 6.3: Depiction the simulation setup for each host in the chain.

6.4 Evaluation data

We want to look at a variety of attack scenarios to highlight the strengths and weaknesses of different SSD approaches. We created three main attack datasets that contain different forms and amounts of evasive behaviour, and a smaller dataset to highlight the influence of different chain lengths.

To present a valuable false positive test, we provide three datasets with benign background traffic. The first contains general real-world traffic, while the second and third contain benign data that bears similar traffic characteristics as the generated attack data.

6.4.1 Stepping-stone data

We create our main datasets using a chain of four stepping-stones S_1, S_2, S_3 , and S_4 . We subdivide into three datasets: We first capture data without transfer delays and chaff perturbations in **dataset BA (baseline attack)**. We then capture data once with added transfer delays with varying δ_D to control delays in **dataset DA (delay attack)**, and once with added chaff perturbations of varying δ_C in **dataset CA (chaff attack)**. Each dataset contains 30.000 connection pairs. We furthermore create a smaller **dataset CL (chain length)** with differing numbers of stepping-stones (1,3,5, and 8 jumps) without transfer delays and chaff perturbations.

6.4.2 Benign data

We include real-world traffic traces, taken from the **CAIDA 2018 Anonymized Internet Traces** dataset [43], as overall background traffic. This data contains traces collected from high-speed monitors on a commercial backbone link, and is often used for research on the characteristics of Internet traffic.

To sufficiently test for false-positive, we also need to include benign traffic that

	Label	Nr. of conn.	Purpose
Attack data	set BA	30 000	Baseline attack data without evasive tactics
	set DA	30 000	Inclusion of delays with varying δ_D
	set CA	30 000	Inclusion of chaff with varying δ_C
	set CL	40 000	Data from chains of different lengths, no evasive tactics
Benign data	CAIDA	60 000	General background data
	SSH	20 000	Background data similar to attack commands
	Multim.	20 000	Background data similar to chaff perturbations

Table 6.1: Summary of different components in our evaluation data.

has similar characteristics to the attack traffic and was generated in a similar network environment. We created a set of interactive SSH-connections that communicate directly between the client and the server without a stepping-stone. We follow the same procedure as described in Section 6.3.2.1.

Since we generate perturbations with multimedia streams characteristics, we additionally want to test for false-positives against actual multimedia stream traffic. For that, we captured traffic from a Nginx-server streaming randomised video to a client.

We merge the three datasets to create our benign background dataset, with the CAIDA part containing 60 000 connection pairs, while the other two each contain 20 000 connection pairs. The amount of SSH traffic and multimedia streams in this setting is inflated from a realistic setting (up to 0.2% of flows for SSH and up 3% for video streaming [145]) to highlight the strengths and drawbacks of SSD methods, which we consider in the evaluation. In Section 6.6.4, we analyse false-positives for each dataset individually. Table 6.1 summarises the different parts in our evaluation data.

6.4.3 Evaluation methodology

To create a fair playing field for the selected SSD methods, we only look at connections that exchange more than 1 500 packets and exclude shorter connections from both the data. The number of packets necessary for detection should ideally be a low possible to enable early detection. The chosen number of 1 500 packets seems like a suitable

minimal limit since all of the selected methods are designed to make successful detection with 300-1 500 packets. Furthermore, there were no connections with less packets in the stepping-stone dataset. True stepping stone connections are rare compared to benign ones, making their detection an imbalanced classification problem. An appropriate evaluation measure for imbalanced data are false positive and false negative rates as well as the *Area-under-ROC-curve* (AUC) for threshold-based methods.

6.5 Selected SSD methods and Implementation

A range of underlying techniques exist for SSD, and we try to include approaches from every area to create an informative overview and highlight strengths and weaknesses. We surveyed publications to create a collection of SSD methods. We started with the publications from surveys [146, 147], and then added impactful recent publications found via Google Scholar (keywords “connection”, “correlation” “stepping-stone”, “detection”, “attack”, “chaff perturbation”). From here, we selected approaches based on the following criteria:

1. The achieved detection and false positive rates claimed by the authors,
2. and whether the model design shows robustness against any evasion tactics as claimed by the authors.
3. We always selected the latest versions if a method has been improved by the authors.

Table 6.2 contains a summary of the included methods. Especially for traditional packet-correlation as well as robust watermarking and anomaly-based methods, there has been little developments since the early 2010s. We labelled each method to make referring to it in the evaluation easier.

PContext, 2011 Yang et al. [148] compare sequences of interarrival times in connection pairs to detect potential stepping-stone behaviour. For that, the contextual distance of a packet is defined as the packet interarrival times around that packet. The authors focus on *Echo*-packets instead of *Send*-packets to resist evasion tactics. The authors evaluate their results with up to 100% chaff ratio with 100% detection rate.

Category	Approach	TP	FP	Robustness	Label
Packet-corr.	Yang, 2011 [148]	100%	0%	jitter/< 80% chaff	PContext
Neural netw.	Nasr, 2018 [16]	90%	0.0002%	small jitter	DeepCorr
	Wu, 2010 [149]	100%	0%	-	WuNeur
RTT-based	Yang, 2015 [140]	not provided		50% chaff	RWalk
	Huang, 2016 [150]	85%	5%	-	Crossover
Anomaly	Crescenzo, 2011 [139]	99%	1%	jitter/chaff	Ano1
	Huang, 2011 [151, 152]	95%	0%	> 25% chaff > 0.2s jitter	Ano2
Watermark	Wang, 2011 [153]	100%	0.5%	< 1.4s jitter	WM

Table 6.2: Summary of included SSD-methods along with the claimed true positive and false positive rates and evasion robustness by the corresponding authors. We added labels to each method for later reference.

WuNeur, 2010 Wu et al. [149] propose a neural network model based on sequences of RTTs, which are fed into a feed-forward network to predict the downstream length of the chain. The network itself only contains one hidden layer and achieves good results only if RTTs are small, i.e. when the stepping-stone chain is completely contained within one LAN-network.

DeepCorr, 2018 Nasr et al. [16] train a deep convolutional neural network to identify connection correlation from the interarrival times and packet sizes in each connection. The trained network is large with over 200 input filters, and consists of three convolutional and three feed-forward layers. On stepping-stones, the authors achieve a 90% detection rate with 0.02% false positives.

RWalk, 2015 Yang et al. [140] combine packet-counting methods and RTT mining methods to improve detection results from [154]. The model resists chaff perturbation by estimating the number of round-trips in a connection via packet-matching and clustering to determine if the connection is being relayed.

C-Over, 2016 Huang et al. [150] use the fact that in long connection chain, the round-trip-time of a packet may be longer than the intervals between two consecutive keystrokes. This will result in cross-overs between request and response, which causes the curve of sorted upstream RTTs to rise more steeply than in a regular connection.

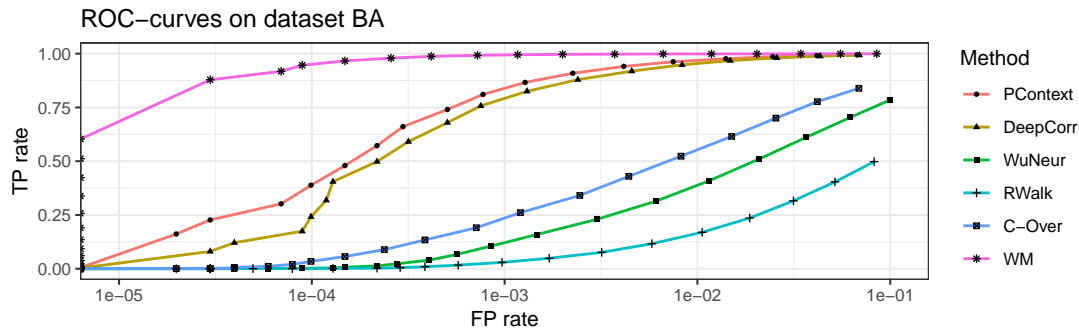


Figure 6.4: ROC-curves for different SSD methods on dataset BA (no evasive tactics). Anomaly-based methods are excluded.

Ano1, 2011 Crescenzo et al. [139] have proposed an anomaly-based methods to detect time delays and chaff perturbations in a selected connection. Packet time-delays are detected if RTTs exceed a threshold, while chaff detection compares the similarity of downstream with upstream sequences. The authors claim detection for chaff ratios 25% or more, and for delays introduced to up to 70% of all packets.

Ano2, 2011/2013 Huang et al. [151, 152] proposed an anomaly-based method to detect chaff and delay perturbations since interarrival times in regular connections tend to follow a Pareto or Lognormal distribution, which chaffed connections supposedly do not. The authors state 95% detection rate at 50% chaff ratio and more while retaining zero false positives using a small set of interactive SSH stepping-stone connections.

WM, 2010 Watermarking typically yields very low false-positives for connection correlation. Wang et al. [153] provide an approach that offers at least some resistance against timing perturbations. The authors assume some limits to an adversary's timing perturbations, such as a bound on the delays. The authors state 100% TP with 0.5% FP with resistance against timing perturbations of up to $1.4s$.

6.6 Results

6.6.1 Data without evasion tactics

First, we look at the detection rates for traffic from stepping-stones that did not use any evasive tactics, i.e. S_1, \dots, S_4 are only forwarding commands and responses. The successful detection of this activity with low false-positives should be the minimum

	PContext	DeepCorr	WuNeur	RWalk	C-Over	WM
AUC	0.998	0.997	0.938	0.853	0.965	0.9998

Table 6.3: AUC-scores for different methods on stepping-stone data without evasive tactics.

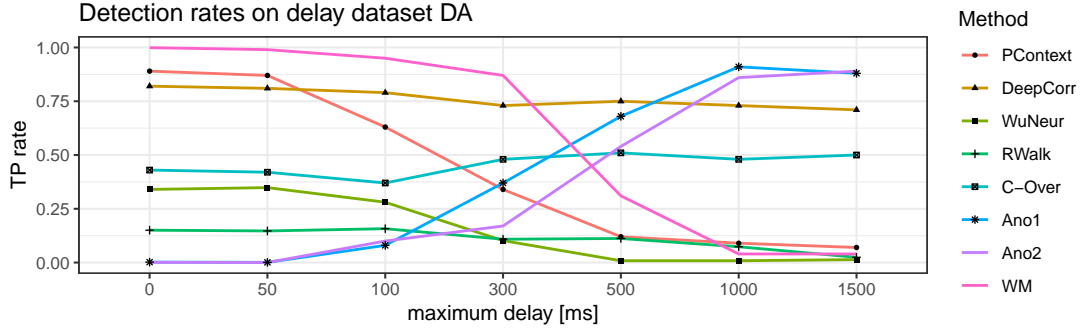


Figure 6.5: Detection rates in dependence of δ_D for different methods on dataset DA with a fixed FP rate of 0.4%.

requirement for any SSD method. Since anomaly-based approaches aim to only detect evasive behaviour, we exclude them from this analysis.

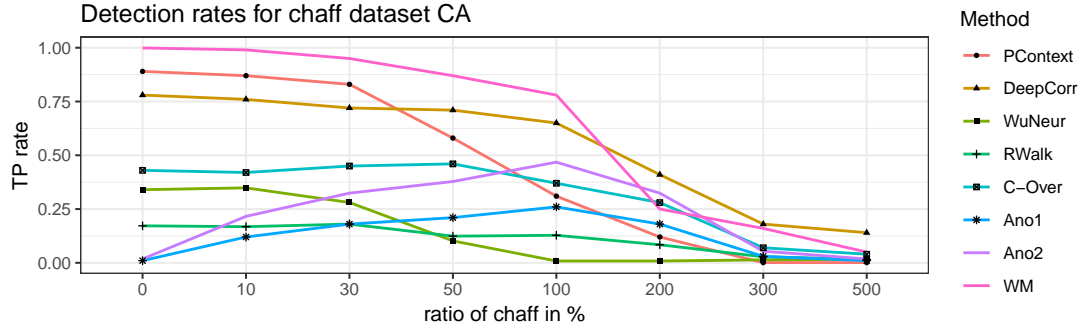
Fig. 6.4 depicts the calculated ROC-curves, which plot the true positive rate against the false positive rate for varying detection thresholds. Table 6.3 depicts the overall AUC-scores.

Unsurprisingly, the watermarking method achieves high detection results with very low false-positives. Both the PContext and DeepCorr models start to yield good detection results of around 80% at a FP rate lower than 0.1%, with the PContext method slightly outpacing the DeepCorr method. RTT-based methods seem to not perform as well compared to the other included methods. Overall, the observed ROC curves seem to be in agreement with the stated detection rates of the selected methods except for RWalk.

6.6.2 Delays

We now consider the effect of transfer delays added by the attacker to packets on the detection rates. For that, we pick detection thresholds for each SSD methods corresponding to a FP rate of 0.4% as most methods are able to achieve at least moderate detection results at this rate. We look at delays added to only to outgoing packets on S_4 , the last stepping stone in the chain. Fig. 6.5 depicts evolution of detection rates in dependence of the maximum delay δ_D .

	PContext	DeepCorr	WuNeur	RWalk	C-Over	Ano1	Ano2	WM
AUC	0.638	0.995	0.613	0.641	0.952	0.997	0.996	0.562

Table 6.4: AUC-scores for SSD methods with added transfer delays at $\delta_D = 1000ms$.Figure 6.6: Detection rates in dependence of δ_C for different methods on dataset CA with a fixed FP rate of 0.4%

As visible, both anomaly-based methods are capable of detecting added delays relatively reliably above a certain threshold. Furthermore, both the detection rates of DeepCorr and the RTT-based C-Over only decrease slightly under the influence of delays. Detection rates for all other methods decrease significantly to the point where no meaningful predictions can be made. This is also reflected by the AUC-scores for traffic with $\delta_D = 1000ms$, given in Table 6.4.

While the WM method is robust against transfer delays up to $\delta_D = 500ms$, this value is smaller than the one claimed by the authors. This might however be a result of the slightly smaller quantisation step size that we used. It is surprising that the PContext method shows only little robustness against transfer delays, which contradicts the authors claims, potentially due to the incorrect assumption that relying on *Echo*-packets are not subject to transfer delays.

6.6.3 Chaff

We now consider the effect of chaff perturbations added by the attacker to individual connections on the detection rates. Again we pick detection thresholds for each SSD methods corresponding to a FP rate of 0.4%.

Chaff packets are added to both the connection between S_3 and S_4 as well as between S_4 and host T as described in Section 6.3.3.2. Fig. 6.6 depicts evolution of detection rates in dependence of the ratio of number of chaff packets to packets from the actual interaction.

	PContext	DeepCorr	WuNeur	RWalk	C-Over	Ano1	Ano2	WM
AUC	0.639	0.886	0.615	0.641	0.589	0.782	0.738	0.839

Table 6.5: AUC-scores for SSD methods with added chaff at 300% ratio.

	PContext	DeepCorr	WuNeur	RWalk	C-Over	Ano1	Ano2	WM
CAIDA	0.36	0.46	0.47	0.67	0.53	0.48	0.35	0.81
SSH	0.53	0.46	0.21	0.28	0.27	0.05	0.02	0.08
multim.	0.11	0.08	0.32	0.04	0.20	0.47	0.63	0.11

Table 6.6: Relative contribution in % of different benign data to the FP rate.

As visible, all methods struggle to detect stepping stones once the chaff packets become the majority of the transferred traffic. This is also evident from the AUC-scores given in Table 6.5. Several approaches claimed to be resistant to chaff perturbations, however prior evaluations were limited chaff ratios below 100% without obvious reason.

It is surprising that the anomaly detection methods do not perform better at detecting chaff perturbations. Chaff in both approaches was however evaluated with different traffic generation distribution and not compared against a background of traffic following a similar generation distribution, which could explain the disagreement between the results we are finding here.

Overall, these results are in disagreement with the "robustness" claims made for four of the selected approaches, namely PContext, RWalk, Ano1, and Ano2.

6.6.4 False positives

Table 6.6 depicts the relative contribution (after adjusting for their weight) at $FP = 0.4\%$ of each of the three benign data types to the overall false positive rate. Most methods have more problems with the heterogeneous nature the CAIDA traces, with only PContext and DeepCorr seeing most false positives in the SSH traffic.

The multimedia traffic is causing most problems for the anomaly-based methods, presumably because it follows a similar distribution as the generated chaff perturbations.

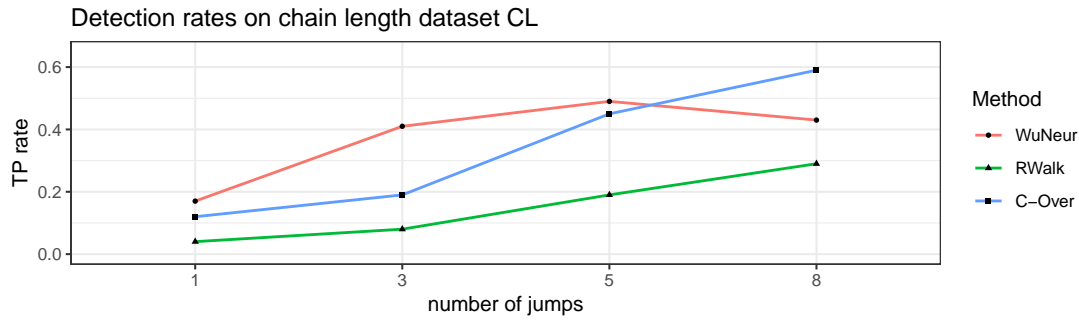


Figure 6.7: Detection rates in dependence of chain length for different methods on dataset CL with a fixed FP rate of 0.4%

	Value	TP deviation from average				
		DeepCorr	WuNeur	RWalk	C-Over	WM
RTT	5ms	−0.2%	+41.3%	−42.3%	−36%	+0.03%
	70ms	−5.6%	−5.8%	+35.1%	+51%	−2.2%
Packet loss	0%	+1.2%	+1.3%	+2.1%	+4.3%	+0.02%
	7%	−9.1%	−1.1%	−3.1%	−7.3%	−9.7%

Table 6.7: Influence of network congestion on detection rates at a fixed FP rate of 0.4%. The given percentages are describing the change of the detection rate under the given congestion setting when compared to the overall average.

6.6.5 Influence of chain length

In this section, we look at the effect of differing chain lengths on the detection rates. We only focus on RTT-based methods here since the other methods should and do not see a significant effect from varying chain lengths. Since RTT-based methods aim to measure the effect of packets travelling via multiple hosts, it is unsurprising that they perform better at detecting longer chains.

Of the RTT-based methods, only C-Over was able to yield consistent detection rates under transfer delays. Interestingly, if the C-Over method is applied to connections between S_3 and S_4 instead of between S_4 and the target, detection rates decrease in the same manner as for other RTT-based methods. This is not surprising as the underlying assumption for robustness for this approach relies on Echo-packets not being delayed.

6.6.6 Influence of network settings

Finally, we look at the effect of different network settings. We only show methods that show significant effects and omitted bandwidth from the evaluation as different values do not seem to have any effect on detection rates¹.

As visible in Table 6.7, the three RTT-based methods show different responses to small/large average round-trip-times. While WuNeur, as expected from prior results, performs better in LAN settings, detection rates of the RWalk and C-Over methods are boosted by larger RTTs. All methods profit from lower packet losses.

6.6.7 Summary

Overall, detection rates on dataset BA are mostly in line with the claimed capabilities except for RWalk, although detection rates are slightly lower than stated by most authors. Delay perturbation increases detection difficulty for most methods, except for Ano1, Ano2, and DeepCorr, which contradicts robustness claims for PContext and to some extent WM. Our inserted chaff perturbations however render detection impossible for all methods examined, which contradicts robustness claims for PContext, Ano1, Ano2, and RWalk, even though the claims were based on lower chaff levels.

As discussed in Section 6.6.5 and 6.6.6, longer chains yield higher detection rates for RTT-based methods while Different network transmission settings seem to have overall little influence on detection rates.

6.7 Conclusion

In this work, we set out to evaluate the state-of-the-art of SSD methods using a comprehensive data generation framework. Our framework simulates realistic stepping-stone behaviour with SSH-tunnels in different settings and varying amounts of evasive perturbation tactics. We released a large dataset that highlights multiple aspects in SSD, and is suitable to train ML-based methods.

Overall, our results show that attackers can reliably evade detection by using the right type and amount of chaff perturbation, which disproves several claims made about the robustness against this evasive tactic. Although to a lesser degree, our implemented delay perturbations still affect detection rates for most methods.

¹For all methods, the detection rate differences (0.7% – 6.2%) were smaller across bandwidths than the overall detection rate errors (2.6% – 6.5%).

Currently, it seems that watermarking methods are most suited to reliably detect simple stepping-stones in real-life deployment. The performance of DeepCorr indicates that deep neural networks show the most potential at detecting attacks that use chaff or delay perturbations if they are trained on suitable data. We find that detection and false-positive rates for RTT-based methods are significantly lower than for other methods.

Chapter 7

Conclusion

This thesis investigated the generation and control of traffic microstructures, their influence on network intrusion detection, and the design and robustness of corresponding models. The presented results both demonstrate how to improve model design and development by considering and examining the difficulties models can experience when encountering specific microstructures, as well as how a model that is designed to capture short-term flow microstructures can raise the state-of-the-art in detecting network access attacks. We also demonstrated how an attacker can perturb microstructures to evade detection in a stepping-stone attack.

In this chapter, we will revisit the original research objectives, discuss some caveats of our work, and suggest future research directions.

7.1 Revisiting research objectives

Considering the presented results in Chapters 3, 4, 5, and 6, we will now discuss how well the research objectives formulated in Sect. 1.3 have been addressed.

Research Objective 1

How well-structured is the space of microstructures observed in the traffic of a machine or a network? To what degree are these microstructures a result of specific computational activities that are of interest for traffic classification and network intrusion detection, and how much are they affected by other external variables?

The results in Chapter 4 provide an examination of factors that shape traffic microstructures, and which ones can be of interest for designing intrusion detection models. The chapter provides experimental evidence for the effect that congestion, compu-

tational load, or background activities can have on packet and connection sequences, or how specific activities such as HTTP-multiplexing or file-syncing or communication failures such as half-open connections can lead to structurally similar traffic that affects the ability of a model to classify correctly. Chapter 5 discusses in part how the length and complexity of connection sequences is determined by specific activities and the amount of activity overlay. It also discusses how frequent different levels of connection sequence complexity are observed in real-world datasets.

Research Objective 2

How can we identify microstructures significant for intrusion detection and train corresponding models? What requirements must a labelled traffic generation framework fulfil to provide realistic data?

Chapter 3 and 5 address this objective to a satisfying degree, albeit not completely. Chapter 3 defines a set of requirements for datasets to enable effective training of machine learning models. These requirements address both the information available on the given data as well as the variability of structures in the data. Chapter 5 provides an example of how to identify a suitable set of microstructures in connection sequences, and how to design a corresponding model to capture them.

Research Objective 3

To what degree can relevant microstructures in network traffic be captured in a model from a training dataset, and how can we achieve this? How can a model adapt to changes of structures in benign traffic?

Chapter 4 provides an examination of how and why two state-of-the-art methods are failing to capture some relevant packet-level microstructures from their given training dataset, and how this impedes their performances. Chapter 5 then examines how our CBAM-model captures flow-level microstructures, what design choices are made to improve the processing of specific structures, and which structures are more difficult to capture. The chapter also provides an examination of how constant these structures remain over time, but does not discuss how to detect and adapt to structural changes. For these, we would require traffic data that introduces structural changes on several hosts, and examine if we can identify characteristics in the way CBAM reacts to these changes compared to anomalies corresponding to attacks.

Research Objective 4

To what degree can access attacks be detected by a model that learns traffic microstructures? What kind of attacks will necessarily show contextual anomalies, and which will not? Can an adversary adapt his attacks to avoid detection?

Chapter 5 provides detection results for access attacks contained in the CICIDS-17 dataset, and discusses reasons why some attacks are detected more effectively than others, and how attackers can potentially evade detection. This is however constrained to the available attack data and not a comprehensive discussion of access attacks in general. Chapter 6 examines a specific type of attack and available detection methods. It also discusses evasive methods extensively and concludes that an attacker can avoid detection in a stepping-stone scenario effectively. In general, there are many future questions to explore regarding the adversarial robustness of detection for different attack types.

7.2 Critical analysis

In this section, we reflect on the limitations of the presented results and how researchers can and cannot benefit from them.

7.2.1 Usage of synthetic data

Chapter 5 presented a novel approach to detect access attacks using network flow data. While we paid attention to include real-world datasets in the longterm evaluation, a lack of available data means that the presented detection rates have been obtained only using synthetic attack data. These attacks were conducted in a laboratory environment, and do not stand in any context to each other. The realism of the delivery of laboratory attacks has been criticised due to the absence of prior or posterior events in relation to the attack, as well as the inflated ratio of attack data. For this reason, the presented detection rates have to be taken with care. The same applies for the methods examined in Chapter 4.

Chapter 6 similarly presents an evaluation of methods on partly real-world benign traffic, but only synthetic attack data that we generated. While similar arguments can be made regarding the realism of the attack traffic, detecting attacks in the wild is generally seen as more challenging than in a laboratory setting. The presented data served as a baseline that a model must be able to fulfil as a minimal requirement. The

conclusion that overall all methods failed to pass this benchmark therefore indicates their unsuitability in a real-world setting.

7.2.2 Operational considerations

Chapter 5 presented some considerations on the computational scalability as well as the maximal detection lag, and Chapter 6 discussed some immediate potential defensive actions following the detection of a stepping stone. However, more work needs to be conducted to confirm that the presented results would be achievable in an operational setting.

7.2.3 Potential incompleteness of microstructure control

Chapter 3 and 4 presented and examined DetGen and its ability to control important factors that shape traffic microstructures. The considered factors were chosen based on domain expertise and the measurement of their respective influence. However, it is difficult to verify that this list is comprehensive and that there are no other important factors that shape traffic microstructures in real-world settings. These might include geographical location and corresponding packet routing or the impact of proxies or firewalls.

For researchers, this means that DetGen is a useful tool to probe models, examine the effect of different factors on model behaviour, or generate sufficient amount and variability of specific traffic types for model training. It is however not necessarily suitable to generate comprehensive datasets that simulate the heterogeneity observed in real-world datasets due to the possible absence of some traffic-shaping influence factors.

7.3 Future work

7.3.1 Future application of DetGen

DetGen currently only generates short-term traffic events. We are investigating how to better simulate causality in connection spawning and other medium-term temporal dependencies, as well as emulate the usage activity of individual scenarios by a user or a network. The modelling and generation of computer network activity has been

investigated extensively, with tools such as DoppelGANger [55] automatically generating realistic network activity streams. A promising way forward would be to enable the import of externally generated activity timelines, and generate the corresponding communication events corresponding to the activity timestamps with DetGen.

An interest of researchers might be to test their NID-model on a variety of network topologies and potentially in real-time. An idea we are investigating is to use DetGen in combination with a virtual network framework such as *Mininet* with virtual software switches, Ethernet links, routers, and firewalls. The idea is to populate the network in a given topology with containers that execute DetGen scenarios with corresponding activity randomisation in real-time according to an activity timeline. This would enable researchers to explore different topologies and place their NID-model within the created network to perform detection in real-time.

Attacks in synthetic NID-datasets so far have been conducted in isolation and do not reflect the chain of attack that a sophisticated intruder might pursue when infiltrating a network. The flexible nature of implemented scenarios in DetGen could be used to execute attack scenarios at different stages and points in a network to simulate the behaviour of a larger attack chain. This could be performed using pre-defined attack tactics to enable different attack paths.

7.3.2 Model and data adaptivity

As applications and internet protocols evolve, so do their corresponding traffic microstructures. We have outlined in Chapter 3 how to adapt data generation to include new traffic types with DetGen's modular design, and in Chapter 5 how to test whether the structures learned by a model remain stable over longer time periods. Additionally, a number of NID-models have been proposed that automatically adapt to slow changes in network characteristics such as port entropies or activity levels [155, 156].

The identification of changes in traffic microstructures, and their distinctions from intrusion attempts remains a challenge. New or updated applications or protocols can alter tokenised packet or flow sequences significantly with novel authentication mechanisms or data retrieval implementations such as HTTP-multiplexing. As these changes are abrupt, their identification as application evolution rather than anomalies caused by intrusion attempts is a complex problem.

7.3.3 Further application of NLP-models and large-scale datasets

Language models based on LSTM-encoders and today's transformer networks such as "BERT" and "GPT3" have demonstrated how well complex structures in tokenised sequences such as written language can be learned when trained on sufficient data. The extensive knowledge of today's language models about structures in the trained data can be used to make accurate probabilistic statements about an input sequence, and is applied in spelling and grammar correction [157], and their design has found application to fraud detection and other anomaly detection applications [158, 159].

The format similarity of packet and flow sequences to sequences of words suggests that success in these areas can translate well to network anomaly detection. With CBAM, we presented a small model that is inspired from NLP and has proven to successfully learn structures in flow sequences. The immense success of current language models such as "BERT" and "GPT3" however stems in part from their large size and the corresponding scale of the datasets they are trained on. Network intrusion detection currently lacks structured large-scale datasets of network traffic to enable similar self-supervised models to thrive in this area. A necessary step to successfully apply state-of-the-art language models to network anomaly detection therefore is the gathering and release of larger and richer datasets that represent and store structures in traffic traces in an efficient and privacy-preserving manner.

7.3.4 Directions for detecting stepping-stones

Our results in Chapter 6 have shown that current methods of detecting interactive stepping-stones fail when encountering evasive chaff perturbations. The ability of an attacker to artificially decorrelate traffic streams suggests that even more sophisticated stream correlation methods could be evaded with sufficient noise. To counter detection evasion, methods that incorporate more contextual features about a network appear more capable of identifying relayed traffic. As the addition of noise produces traffic with a substantial volume, evaluating whether voluminous connections between the respective source and target device in the network should be seen as anomalous. Graph-based methods discussed in Section 6.2 go in a similar direction and could potentially be combined with weak correlation indicators.

7.4 Final outlook

In the end we believe that the study of traffic microstructures and their influence on model behaviour plays a significant role in closing the semantic gap that still exists between the design of ML-based NID techniques and the challenges we encounter when dealing with real-world traffic. Despite 30 years of research, several problems that have been described in this thesis still prevent the overwhelming success that machine and deep learning has experienced in other domains to transfer into better defence against cyber attacks. However, we believe that these problems can be addressed in the future as they are not inherent to the detection of cyber attacks, but are caused by the limited amount of empirical knowledge and data resources available to researchers at the moment.

Bibliography

- [1] H. Clausen, R. Flood, and D. Aspinall, “Traffic generation using Containerization for Machine Learning,” in *Proceedings of the ACSAC Dynamic and Novel Advances in Machine Learning and Intelligent Cyber Security Workshop*. ACM, 2019.
- [2] H. Clausen and D. Aspinall, “Examining traffic microstructures to improve model development,” in *2021 WTMIC at IEEE Security and Privacy Workshops (SPW)*. IEEE, 2021.
- [3] H. Clausen, R. Flood, and D. Aspinall, “Controlling network traffic microstructures for machine-learning model probing,” 2021, manuscript submitted for publication.
- [4] H. Clausen, G. Grov, M. Sabate, and D. Aspinall, “Better Anomaly Detection for Access Attacks Using Deep Bidirectional LSTMs,” in *Machine Learning for Networking: Third International Conference, MLN 2020, Paris, France, November 24–26, 2020, Revised Selected Papers 3*. Springer International Publishing, 2021, pp. 1–18.
- [5] H. Clausen, G. Grov, and D. Aspinall, “Cbam: A contextual model for network anomaly detection,” *Computers*, vol. 10, no. 6, p. 79, 2021.
- [6] H. Clausen, M. S. Gibson, and D. Aspinall, “Evading stepping-stone detection with enough chaff,” in *International Conference on Network and System Security*. Springer, 2020, pp. 431–446.
- [7] “NCSC CAF guidance,” <https://www.ncsc.gov.uk/collection/caf/caf-principles-and-guidance/c-2-proactive-security-event-discovery>, UK National Cyber Security Centre, Tech. Rep., 2019, accessed on 25 August 2021.

- [8] M. Roesch *et al.*, “Snort: Lightweight intrusion detection for networks.” in *Lisa*, vol. 99, no. 1, 1999, pp. 229–238.
- [9] C. Nachreiner, “Internet Security Report - Q1 2021,” <https://www.watchguard.com/wgrd-resource-center/security-report-q1-2021>, WatchGuard Technologies, Tech. Rep., 2021, accessed on 25 August 2021.
- [10] “Catching Mimikatz’ behavior with anomaly detection,” <https://www.darktrace.com/en/blog/catching-mimikatz-behavior-with-anomaly-detection/>, Darktrace, Tech. Rep., 2019, accessed on 22 July 2021.
- [11] “How the CrowdStrike ML-based Engine Defends Against Unknown Malware,” <https://www.crowdstrike.com/resources/white-papers/rise-machine-learning-ml-cybersecurity/>, CrowdStrike, Tech. Rep., 2020, accessed on 22 July 2021.
- [12] UK National Cyber Security Centre Guidance, “Intelligent security tools,” <https://www.ncsc.gov.uk/collection/intelligent-security-tools>, Tech. Rep., 2019, accessed 23 August 2021.
- [13] A. Nisioti, A. Mylonas, P. D. Yoo, and V. Katos, “From Intrusion Detection to Attacker Attribution: A Comprehensive Survey of Unsupervised Methods,” *IEEE Communications Surveys & Tutorials*, 2018.
- [14] R. Sommer and V. Paxson, “Outside the Closed World: On Using Machine Learning for Network Intrusion Detection,” in *2010 IEEE Symposium on Security and Privacy*, May 2010, pp. 305–316.
- [15] Y. Mirsky, T. Doitshman, Y. Elovici, and A. Shabtai, “Kitsune: an ensemble of autoencoders for online network intrusion detection,” *arXiv preprint arXiv:1802.09089*, 2018.
- [16] M. Nasr, A. Bahramali, and A. Houmansadr, “Deepcorr: Strong flow correlation attacks on tor using deep learning,” in *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*, 2018, pp. 1962–1976.
- [17] R. Harang, “Bridging the semantic gap: Human factors in anomaly-based intrusion detection systems,” in *Network Science and Cybersecurity*. Springer, 2014, pp. 15–37.

- [18] M. R. Smith, N. T. Johnson, J. B. Ingram, A. J. Carbajal, R. Ramyaa, E. Domschot, C. C. Lamb, S. J. Verzi, and W. P. Kegelmeyer, “Mind the Gap: On Bridging the Semantic Gap between Machine Learning and Information Security,” *arXiv preprint arXiv:2005.01800*, 2020.
- [19] N. Brownlee, C. Mills, and G. Ruth, “Traffic flow measurement: Architecture,” Tech. Rep., 1999.
- [20] J. P. Anderson, “Computer security threat monitoring and surveillance,” *Technical Report, James P. Anderson Company*, 1980.
- [21] V. Chandola, A. Banerjee, and V. Kumar, “Anomaly Detection: A Survey,” *ACM Comput. Surv.*, vol. 41, no. 3, pp. 15:1–15:58, Jul. 2009.
- [22] A. v. d. Oord, Y. Li, and O. Vinyals, “Representation learning with contrastive predictive coding,” *arXiv preprint arXiv:1807.03748*, 2018.
- [23] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, “Albert: A lite bert for self-supervised learning of language representations,” *arXiv preprint arXiv:1909.11942*, 2019.
- [24] D. E. Denning, “An intrusion-detection model,” *IEEE Transactions on software engineering*, no. 2, pp. 222–232, 1987.
- [25] M. M. Sebring, “Expert systems in intrusion detection: A case study,” in *Proc. 11th National Computer Security Conference, Baltimore, Maryland, Oct. 1988*, 1988, pp. 74–81.
- [26] S. E. Smaha *et al.*, “Haystack: An intrusion detection system,” in *Fourth Aerospace Computer Security Applications Conference*, vol. 44. Orlando, FL, USA, 1988.
- [27] A. Wagner, T. Dübendorfer, B. Plattner, and R. Hiestand, “Experiences with worm propagation simulations,” in *Proceedings of the 2003 ACM workshop on Rapid Malcode*, 2003, pp. 34–41.
- [28] A. Wagner and B. Plattner, “Entropy based worm and anomaly detection in fast IP networks,” in *Enabling Technologies: Infrastructure for Collaborative Enterprise, 2005. 14th IEEE International Workshops on.* IEEE, 2005, pp. 172–177.

- [29] A. Lakhina, M. Crovella, and C. Diot, “Diagnosing Network-wide Traffic Anomalies,” in *Proceedings of the 2004 Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications*, ser. SIGCOMM ’04. New York, NY, USA: ACM, 2004, pp. 219–230, 01230.
- [30] —, “Characterization of Network-wide Anomalies in Traffic Flows,” in *Proceedings of the 4th ACM SIGCOMM Conference on Internet Measurement*, ser. IMC ’04. New York, NY, USA: ACM, 2004, pp. 201–206, 00439.
- [31] C. Krügel, T. Toth, and E. Kirda, “Service specific anomaly detection for network intrusion detection,” in *Proceedings of the 2002 ACM symposium on Applied computing*. ACM, 2002, pp. 201–208.
- [32] M. V. Mahoney and P. K. Chan, “Learning nonstationary models of normal network traffic for detecting novel attacks,” in *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2002, pp. 376–385.
- [33] MIT Lincoln Labs, “1998 DARPA Intrusion Detection Evaluation,” <https://www.ll.mit.edu/r-d/datasets/1998-darpa-intrusion-detection-evaluation-dataset>, 1998, accessed: 2021-08-11.
- [34] University of California, Irvine, “KDD Cup 1999,” <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>, 1999, accessed: 2021-08-11.
- [35] M. V. Mahoney and P. K. Chan, “An analysis of the 1999 DARPA/Lincoln Laboratory evaluation data for network anomaly detection,” in *International Workshop on Recent Advances in Intrusion Detection*. Springer, 2003, pp. 220–237.
- [36] T. Shon and J. Moon, “A hybrid machine learning approach to network anomaly detection,” *Information Sciences*, vol. 177, no. 18, pp. 3799–3821, 2007.
- [37] N. Shone, T. N. Ngoc, V. D. Phai, and Q. Shi, “A deep learning approach to network intrusion detection,” *IEEE transactions on emerging topics in computational intelligence*, vol. 2, no. 1, pp. 41–50, 2018.
- [38] A. D. Kent, “Cybersecurity Data Sources for Dynamic Network Research,” in *Dynamic Networks in Cybersecurity*. Imperial College Press, Jun. 2015.

- [39] ———, “Comprehensive, Multi-Source Cyber-Security Events,” Los Alamos National Laboratory, 2015.
- [40] M. J. M. Turcotte, A. D. Kent, and C. Hash, “Unified Host and Network Data Set,” *ArXiv e-prints*, 2017.
- [41] G. Maciá-Fernández, J. Camacho, R. Magán-Carrión, P. García-Teodoro, and R. Therón, “UGR ‘16: A new dataset for the evaluation of cyclostationarity-based network IDSs,” *Computers & Security*, vol. 73, pp. 411–424, 2018.
- [42] R. Damasevicius, A. Venckauskas, S. Grigaliunas, J. Toldinas, N. Morkevicius, T. Aleliunas, and P. Smuikys, “LITNET-2020: An annotated real-world network flow dataset for network intrusion detection,” *Electronics*, vol. 9, no. 5, p. 800, 2020.
- [43] “The CAIDA UCSD Anonymized Internet Traces 2018,” https://www.caida.org/data/passive/passive_2018_dataset.xml, 2018, accessed on 22 July 2021.
- [44] R. Fontugne, P. Borgnat, P. Abry, and K. Fukuda, “MAWILab: Combining Diverse Anomaly Detectors for Automated Anomaly Labeling and Performance Benchmarking,” in *ACM CoNEXT ’10*, Philadelphia, PA, December 2010.
- [45] UNIBS, “Data Sharing,” <http://netweb.ing.unibs.it/~ntw/tools/traces/>, 2009, accessed on 05 Nov. 2018.
- [46] A. Gharib, I. Sharafaldin, A. H. Lashkari, and A. A. Ghorbani, “An evaluation framework for intrusion detection dataset,” in *Information Science and Security (ICISS), 2016 International Conference on*. IEEE, 2016, pp. 1–6.
- [47] I. Sharafaldin, A. Gharib, A. H. Lashkari, and A. A. Ghorbani, “Towards a reliable intrusion detection benchmark dataset,” *Software Networking*, vol. 2018, no. 1, pp. 177–200, 2018.
- [48] N. Moustafa and J. Slay, “UNSW-NB15: a comprehensive data set for network intrusion detection systems (UNSW-NB15 network data set),” in *2015 Military Communications and Information Systems Conference (MilCIS)*, Nov. 2015, pp. 1–6.
- [49] R. P. Lippmann, D. J. Fried, I. Graf, J. W. Haines, K. R. Kendall, D. McClung, D. Weber, S. E. Webster, D. Wyschogrod, R. K. Cunningham *et al.*, “Evaluating

- intrusion detection systems: The 1998 DARPA off-line intrusion detection evaluation,” in *DARPA Information Survivability Conference and Exposition, 2000. DISCEX’00. Proceedings*, vol. 2. IEEE, 2000, pp. 12–26.
- [50] “The CTU-13 Dataset. A Labeled Dataset with Botnet, Normal and Background traffic.” <http://mcfp.weebly.com/the-ctu-13-dataset.html>, accessed on 22 July 2021.
- [51] S. Garcia, M. Grill, J. Stiborek, and A. Zunino, “An empirical comparison of botnet detection methods,” *computers & security*, vol. 45, pp. 100–123, 2014.
- [52] C.-H. Hsu and U. Kremer, “Iperf: A framework for automatic construction of performance prediction models,” in *Workshop on Profile and Feedback-Directed Compilation (PFDC), Paris, France*. Citeseer, 1998.
- [53] H. Haas, *Netsniff-ng toolkit*, <https://web.archive.org/web/20160908021235/http://netsniff-ng.org/>, 2010, accessed: 2021-08-24.
- [54] Keysight, *PerfectStorm 100GE High-Performance Application and Security Load Modules*, <https://www.keysight.com/de/de/assets/3120-1225/data-sheets/>, accessed: 2021-08-24.
- [55] Z. Lin, A. Jain, C. Wang, G. Fanti, and V. Sekar, “Generating High-fidelity, Synthetic Time Series Datasets with DoppelGANger,” *arXiv preprint arXiv:1909.13403*, 2019.
- [56] A. Cheng, “Pac-gan: Packet generation of network traffic using generative adversarial networks,” in *2019 IEEE 10th Annual Information Technology, Electronics and Mobile Communication Conference (IEMCON)*. IEEE, 2019, pp. 0728–0734.
- [57] Q. Cheng, S. Zhou, Y. Shen, D. Kong, and C. Wu, “Packet-Level Adversarial Network Traffic Crafting using Sequence Generative Adversarial Networks,” *arXiv preprint arXiv:2103.04794*, 2021.
- [58] Y. Mirsky, A. Shabtai, L. Rokach, B. Shapira, and Y. Elovici, “Sherlock vs moriarty: A smartphone dataset for cybersecurity research,” in *Proceedings of the 2016 ACM workshop on Artificial intelligence and security*. ACM, 2016, pp. 1–12.

- [59] R. Fujdiak, V. Uher, P. Mlynek, P. Blazek, J. Slacik, V. Sedlacek, J. Misurec, M. Volkova, and P. Chmelar, "IP Traffic Generator Using Container Virtualization Technology," in *2018 10th International Congress on Ultra Modern Telecommunications and Control Systems and Workshops (ICUMT)*. IEEE, 2018, pp. 1–6.
- [60] D. Inc., "Docker," <https://www.docker.com/>, accessed: 2019-08-11.
- [61] K. Kolyskin, "Virtualization in linux," *White paper, OpenVZ*, vol. 3, p. 39, 2006.
- [62] D. Inc., *Docker Hub*, <https://hub.docker.com/>, accessed: 2019-08-11.
- [63] ———, *Docker Documentation*, <https://docs.docker.com/>, accessed: 2021-08-11.
- [64] S. Hemminger *et al.*, "Network emulation with NetEm," in *Linux conf au*, 2005, pp. 18–23.
- [65] A. Sperotto, R. Sadre, F. Van Vliet, and A. Pras, "A labeled data set for flow-based intrusion detection," in *International Workshop on IP Operations and Management*. Springer, 2009, pp. 39–50.
- [66] K. Allix, T. F. D. A. Bissyande, J. Klein, and Y. Le Traon, "Machine Learning-Based Malware Detection for Android Applications: History Matters!" University of Luxembourg, SnT, Tech. Rep., 2014.
- [67] A. Shiravi, H. Shiravi, M. Tavallaee, and A. A. Ghorbani, "Toward developing a systematic approach to generate benchmark datasets for intrusion detection," *computers & security*, vol. 31, no. 3, pp. 357–374, 2012.
- [68] M. Gates and A. Warshavsky, *Iperf Man Page*, <https://linux.die.net/man/1/iperf>, accessed: 2019-08-11.
- [69] V. Jacobson, C. Leres, and S. McCanne, "The tcpdump manual page," *Lawrence Berkeley Laboratory, Berkeley, CA*, vol. 143, 1989.
- [70] D. Inc., *Best Practises for Writing Docker Images*, https://docs.docker.com/develop/develop-images/dockerfile_best-practices/, accessed: 2021-08-11.
- [71] S. Renfro, *Mergecap Man Page*, <https://www.wireshark.org/docs/man-pages/mergcap.html>, accessed: 2019-08-11.

- [72] R. Sharpe, *Editcap Man Page*, <https://www.wireshark.org/docs/man-pages/editcap.html>, accessed: 2019-08-11.
- [73] F. J. Massey Jr, “The Kolmogorov-Smirnov test for goodness of fit,” *Journal of the American statistical Association*, vol. 46, no. 253, pp. 68–78, 1951.
- [74] M. A. Arfeen, K. Pawlikowski, D. McNickle, and A. Willig, “The role of the weibull distribution in internet traffic modeling,” in *Proceedings of the 2013 25th International Teletraffic Congress (ITC)*. IEEE, 2013, pp. 1–8.
- [75] V. Paxson and S. Floyd, “Wide area traffic: the failure of Poisson modeling,” *IEEE/ACM Transactions on networking*, vol. 3, no. 3, pp. 226–244, 1995.
- [76] S. Zander, T. Nguyen, and G. Armitage, “Automated traffic classification and application identification using machine learning,” in *The IEEE Conference on Local Computer Networks 30th Anniversary (LCN’05) I*. IEEE, 2005, pp. 250–257.
- [77] T. T. Nguyen and G. J. Armitage, “A survey of techniques for internet traffic classification using machine learning,” *IEEE Communications Surveys and Tutorials*, vol. 10, no. 1-4, pp. 56–76, 2008.
- [78] M. Jaber, R. G. Cascella, and C. Barakat, “Can we trust the inter-packet time for traffic classification?” in *2011 IEEE International Conference on Communications (ICC)*. IEEE, 2011, pp. 1–5.
- [79] L. Bernaille, R. Teixeira, I. Akodkenou, A. Soule, and K. Salamatian, “Traffic classification on the fly,” *ACM SIGCOMM Computer Communication Review*, vol. 36, no. 2, pp. 23–26, 2006.
- [80] J. Wexler, M. Pushkarna, T. Bolukbasi, M. Wattenberg, F. Viégas, and J. Wilson, “The what-if tool: Interactive probing of machine learning models,” *IEEE transactions on visualization and computer graphics*, vol. 26, no. 1, pp. 56–65, 2019.
- [81] S. R. Livingstone and F. A. Russo, “The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS): A dynamic, multimodal set of facial and vocal expressions in North American English,” *PloS one*, vol. 13, no. 5, p. e0196391, 2018.

- [82] A. Haque, M. Guo, P. Verma, and L. Fei-Fei, "Audio-linguistic embeddings for spoken sentences," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 7355–7359.
- [83] H. Kamper, Y. Matusevych, and S. Goldwater, "Multilingual acoustic word embedding models for processing zero-resource languages," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 6414–6418.
- [84] I. Sharafaldin, A. H. Lashkari, and A. A. Ghorbani, "Toward generating a new intrusion detection dataset and intrusion traffic characterization," in *ICISSP*, 2018, pp. 108–116.
- [85] N. Moustafa and J. Slay, "UNSW-NB15: a comprehensive data set for network intrusion detection systems," in *2015 military communications and information systems conference (MilCIS)*. IEEE, 2015, pp. 1–6.
- [86] M. Ring, S. Wunderlich, D. Scheuring, D. Landes, and A. Hotho, "A survey of network-based intrusion detection data sets," *Computers & Security*, vol. 86, pp. 147–167, 2019.
- [87] R.-H. Hwang, M.-C. Peng, V.-L. Nguyen, and Y.-L. Chang, "An LSTM-Based Deep Learning Approach for Classifying Malicious Traffic at the Packet Level," *Applied Sciences*, vol. 9, no. 16, p. 3414, 2019.
- [88] H. Tahaei, F. Afifi, A. Asemi, F. Zaki, and N. B. Anuar, "The rise of traffic classification in IoT networks: A survey," *Journal of Network and Computer Applications*, vol. 154, p. 102538, 2020.
- [89] P. Mell, V. Hu, R. Lippmann, J. Haines, and M. Zissman, "An overview of issues in testing intrusion detection systems," 2003.
- [90] X. Hou, L. Shen, K. Sun, and G. Qiu, "Deep feature consistent variational autoencoder," in *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2017, pp. 1133–1141.
- [91] T.-F. Yen, X. Huang, F. Monroe, and M. K. Reiter, "Browser fingerprinting from coarse traffic summaries: Techniques and implications," in *International*

- Conference on Detection of Intrusions and Malware, and Vulnerability Assessment*. Springer, 2009, pp. 157–175.
- [92] T. Stöber, M. Frank, J. Schmitt, and I. Martinovic, “Who do you sync you are? smartphone fingerprinting via application behaviour,” in *Proceedings of the sixth ACM conference on Security and privacy in wireless and mobile networks*, 2013, pp. 7–12.
- [93] R. Marx, J. Herbots, W. Lamotte, and P. Quax, “Same Standards, Different Decisions: A Study of QUIC and HTTP/3 Implementation Diversity,” in *Proceedings of the Workshop on the Evolution, Performance, and Interoperability of QUIC*, 2020, pp. 14–20.
- [94] C. Fricker, P. Robert, J. Roberts, and N. Sbihi, “Impact of traffic mix on caching performance in a content-centric network,” in *2012 Proceedings IEEE INFOCOM Workshops*. IEEE, 2012, pp. 310–315.
- [95] F. J. Aparicio-Navarro, J. A. Chambers, K. Kyriakopoulos, Y. Gong, and D. Parish, “Using the pattern-of-life in networks to improve the effectiveness of intrusion detection systems,” in *2017 IEEE International Conference on Communications (ICC)*. IEEE, 2017, pp. 1–7.
- [96] B. J. Radford, L. M. Apolonio, A. J. Trias, and J. A. Simpson, “Network traffic anomaly detection using recurrent neural networks,” *arXiv preprint arXiv:1803.10769*, 2018.
- [97] L. Sequeira, J. Fernández-Navajas, L. Casadesus, J. Saldana, I. Quintana, and J. Ruiz-Mas, “The influence of the buffer size in packet loss for competing multimedia and bursty traffic,” in *2013 International Symposium on Performance Evaluation of Computer and Telecommunication Systems (SPECTS)*. IEEE, 2013, pp. 134–141.
- [98] A. Jurgelionis, J.-P. Laulajainen, M. Hirvonen, and A. I. Wang, “An empirical study of NetEm network emulation functionalities,” in *2011 Proceedings of 20th International Conference on Computer Communications and Networks (ICCCN)*. IEEE, 2011, pp. 1–6.
- [99] 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, pp. 253–262.

- [100] Y. Aun, S. Manickam, and S. Karuppayah, "A review on features' robustness in high diversity mobile traffic classifications," *International journal of communication networks and information security*, vol. 9, no. 2, p. 294, 2017.
- [101] "M-Trends 2015: A view from the front lines," <https://www2.fireeye.com/rs/fireeye/images/rpt-m-trends-2015.pdf>, Mandiant, Tech. Rep., 2015, accessed on 22 July 2021.
- [102] L. Bontemps, V. L. Cao, J. McDermott, and N.-A. Le-Khac, "Collective Anomaly Detection Based on Long Short-Term Memory Recurrent Neural Networks," in *Future Data and Security Engineering*. Springer, 2016, pp. 141–152.
- [103] J. Kim, J. Kim, H. L. T. Thu, and H. Kim, "Long short term memory recurrent neural network classifier for intrusion detection," in *2016 International Conference on Platform Technology and Service (PlatCon)*. IEEE, 2016, pp. 1–5.
- [104] A. Özgür and H. Erdem, "A review of KDD99 dataset usage in intrusion detection and machine learning between 2010 and 2015," *PeerJ Preprints*, vol. 4, p. e1954v1, 2016.
- [105] M. Tavallaei, E. Bagheri, W. Lu, and A. A. Ghorbani, "A detailed analysis of the KDD CUP 99 data set," in *2009 IEEE Symposium on Computational Intelligence for Security and Defense Applications*. IEEE, 2009, pp. 1–6.
- [106] E. Koutrouki, "Mitigating concept drift in data mining applications for intrusion detection systems," *arXiv preprint arXiv:1010.4784*, 2018.
- [107] P. Rubin-Delanchy, D. J. Lawson, M. J. Turcotte, N. Heard, and N. M. Adams, "Three statistical approaches to sessionizing network flow data," in *2014 IEEE Joint Intelligence and Security Informatics Conference*. IEEE, 2014, pp. 244–247.
- [108] W. Chen, D. Grangier, and M. Auli, "Strategies for training large vocabulary neural language models," *arXiv preprint arXiv:1512.04906*, 2015.
- [109] W. Yin, K. Kann, M. Yu, and H. Schütze, "Comparative study of CNN and RNN for natural language processing," *arXiv preprint arXiv:1702.01923*, 2017.

- [110] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *The journal of machine learning research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [111] M. Ahmed, A. N. Mahmood, and J. Hu, "A survey of network anomaly detection techniques," *Journal of Network and Computer Applications*, vol. 60, pp. 19–31, 2016.
- [112] J. W. Ulvila and J. E. Gaffney Jr, "Evaluation of intrusion detection systems," *Journal of Research of the National Institute of Standards and Technology*, vol. 108, no. 6, p. 453, 2003.
- [113] A. Milenkoski, M. Vieira, S. Kounev, A. Avritzer, and B. D. Payne, "Evaluating computer intrusion detection systems: A survey of common practices," *ACM Computing Surveys (CSUR)*, vol. 48, no. 1, pp. 1–41, 2015.
- [114] L. Liu, P. Wang, J. Lin, and L. Liu, "Intrusion Detection of Imbalanced Network Traffic Based on Machine Learning and Deep Learning," *IEEE Access*, 2020.
- [115] P. Casas, J. Mazel, and P. Owezarski, "Unsupervised network intrusion detection systems: Detecting the unknown without knowledge," *Computer Communications*, vol. 35, no. 7, pp. 772–783, 2012.
- [116] A. Javaid, Q. Niyaz, W. Sun, and M. Alam, "A deep learning approach for network intrusion detection system," in *Proceedings of the 9th EAI International Conference on Bio-inspired Information and Communications Technologies (formerly BIONETICS)*. ICST (Institute for Computer Sciences, Social-Informatics and Telecommunications Engineering), 2016, pp. 21–26.
- [117] N. Ye *et al.*, "A markov chain model of temporal behavior for anomaly detection," in *Proceedings of the 2000 IEEE Systems, Man, and Cybernetics Information Assurance and Security Workshop*, vol. 166. West Point, NY, 2000, p. 169.
- [118] G. Pellegrino, Q. Lin, C. Hammerschmidt, and S. Verwer, "Learning behavioral fingerprints from netflows using timed automata," in *Integrated Network and Service Management (IM), 2017 IFIP/IEEE Symposium on*. IEEE, 2017, pp. 308–316.

- [119] G. Grov, M. Sabate, W. Chen, and D. Aspinall, “Towards Intelligible Robust Anomaly Detection by Learning Interpretable Behavioural Models,” *NISK Journal*, vol. 12, 2019.
- [120] Y. Yu, G. Liu, H. Yan, H. Li, and H. Guan, “Attention-Based Bi-LSTM Model for Anomalous HTTP Traffic Detection,” in *2018 15th International Conference on Service Systems and Service Management (ICSSSM)*. IEEE, 2018, pp. 1–6.
- [121] Y. Song, A. D. Keromytis, and S. Swolfo, “Spectrogram: A mixture-of-markov-chains model for anomaly detection in web traffic,” 2009.
- [122] T. Wakui, T. Kondo, and F. Teraoka, “GAMPAL: Anomaly Detection for Internet Backbone Traffic by Flow Prediction with LSTM-RNN,” in *International Conference on Machine Learning for Networking*. Springer, 2019, pp. 196–211.
- [123] M. A. Khan, “HCRNNIDS: Hybrid Convolutional Recurrent Neural Network-Based Network Intrusion Detection System,” *Processes*, vol. 9, no. 5, p. 834, 2021.
- [124] Y. Zhong, W. Chen, Z. Wang, Y. Chen, K. Wang, Y. Li, X. Yin, X. Shi, J. Yang, and K. Li, “HELAD: A novel network anomaly detection model based on heterogeneous ensemble learning,” *Computer Networks*, vol. 169, p. 107049, 2020.
- [125] X. Zhou, Y. Hu, W. Liang, J. Ma, and Q. Jin, “Variational LSTM enhanced anomaly detection for industrial big data,” *IEEE Transactions on Industrial Informatics*, vol. 17, no. 5, pp. 3469–3477, 2020.
- [126] D. S. Berman, A. L. Buczak, J. S. Chavis, and C. L. Corbett, “A survey of deep learning methods for cyber security,” *Information*, vol. 10, no. 4, p. 122, 2019.
- [127] Y. Shen, E. Mariconti, P. A. Vervier, and G. Stringhini, “Tiresias: Predicting security events through deep learning,” in *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*. ACM, 2018, pp. 592–605.
- [128] M. Du, F. Li, G. Zheng, and V. Srikumar, “Deeplog: Anomaly detection and diagnosis from system logs through deep learning,” in *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*. ACM, 2017, pp. 1285–1298.

- [129] X. Han, T. Pasquier, A. Bates, J. Mickens, and M. Seltzer, “Unicorn: Runtime provenance-based detector for advanced persistent threats,” *arXiv preprint arXiv:2001.01525*, 2020.
- [130] C. Tankard, “Advanced persistent threats and how to monitor and deter them,” *Network security*, vol. 2011, no. 8, pp. 16–19, 2011.
- [131] “McAfee Technical Report on Night Dragon Operation,” https://www.mcafee.com/wp-content/uploads/2011/02/McAfee_NightDragon_wp_draft_to_customersv1-1.pdf, McAfee, Tech. Rep., 2015, accessed on 22 July 2021.
- [132] R. M. Lee, M. J. Assante, and T. Conway, “Analysis of the Cyber Attack on the Ukrainian Power Grid,” https://www.nerc.com/pa/CI/ESISAC/Documents/E-ISAC_SANS_Ukraine_DUC_18Mar2016.pdf, E-ISAC, Tech. Rep., 2016, accessed on 22 July 2021.
- [133] L. Ayala, “Active Medical Device Cyber-Attacks,” in *Cybersecurity for Hospitals and Healthcare Facilities*. Springer, 2016, pp. 19–37.
- [134] E. ENISA, “Baseline Security Recommendations for IoT in the context of Critical Information Infrastructures,” 2017.
- [135] S. Staniford-Chen and L. T. Heberlein, “Holding intruders accountable on the internet,” in *Proceedings 1995 IEEE Symposium on Security and Privacy*. IEEE, 1995, pp. 39–49.
- [136] G. Fraser, “Tunneling, Pivoting, and WebApplication PenetrationTesting,” <https://www.sans.org/reading-room/whitepapers/testing/tunneling-pivoting-web-application-penetration-testing-36117>, SANS, Tech. Rep., 2015, accessed on 22 July 2021.
- [137] D. L. Donoho, A. G. Flesia, U. Shankar, V. Paxson, J. Coit, and S. Staniford, “Multiscale stepping-stone detection: Detecting pairs of jittered interactive streams by exploiting maximum tolerable delay,” in *International Workshop on Recent Advances in Intrusion Detection*. Springer, 2002, pp. 17–35.
- [138] J. Xin, L. Zhang, B. Aswegan, J. Dickerson, T. Daniels, and Y. Guan, “A testbed for evaluation and analysis of stepping stone attack attribution techniques,”

- in *2nd International Conference on Testbeds and Research Infrastructures for the Development of Networks and Communities, 2006. TRIDENTCOM 2006*. IEEE, 2006, pp. 9–pp.
- [139] G. Di Crescenzo, A. Ghosh, A. Kampasi, R. Talpade, and Y. Zhang, “Detecting Anomalies in Active Insider Stepping Stone Attacks.” *JoWUA*, vol. 2, no. 1, pp. 103–120, 2011.
 - [140] J. Yang and Y. Zhang, “RTT-based Random Walk Approach to Detect Stepping-Stone Intrusion,” in *2015 IEEE 29th International Conference on Advanced Information Networking and Applications*. IEEE, 2015, pp. 558–563.
 - [141] B. Coskun and N. Memon, “Efficient detection of delay-constrained relay nodes,” in *Twenty-Third Annual Computer Security Applications Conference (ACSAC 2007)*. IEEE, 2007, pp. 353–362.
 - [142] G. Ho, M. Dhiman, D. Akhawe, V. Paxson, S. Savage, G. M. Voelker, and D. Wagner, “Hopper: Modeling and detecting lateral movement,” in *30th USENIX Security Symposium (USENIX Security 21)*, 2021, pp. 3093–3110.
 - [143] G. Apruzzese, F. Pierazzi, M. Colajanni, and M. Marchetti, “Detection and threat prioritization of pivoting attacks in large networks,” *IEEE Transactions on Emerging Topics in Computing*, 2017.
 - [144] J. D. Padhye, K. Kothari, M. Venkateshaiah, and M. Wright, “Evading stepping-stone detection under the cloak of streaming media with SNEAK,” *Computer Networks*, vol. 54, no. 13, pp. 2310–2325, 2010.
 - [145] P. Velan, J. Medková, T. Jirsík, and P. Čeleda, “Network traffic characterisation using flow-based statistics,” in *NOMS 2016-2016 IEEE/IFIP Network Operations and Management Symposium*. IEEE, 2016, pp. 907–912.
 - [146] R. Shullich, J. Chu, P. Ji, and W. Chen, “A survey of research in stepping-stone detection,” ” *International Journal of Electronic Commerce Studies*”, vol. 2, no. 2, pp. 103–126, 2011.
 - [147] L. Wang and J. Yang, “A research survey in stepping-stone intrusion detection,” *EURASIP Journal on Wireless Communications and Networking*, vol. 2018, no. 1, p. 276, 2018.

- [148] J. Yang and D. Woolbright, “Correlating TCP/IP Packet contexts to detect stepping-stone intrusion,” *Computers & Security*, vol. 30, no. 6-7, pp. 538–546, 2011.
- [149] H.-C. Wu and S.-H. S. Huang, “Neural networks-based detection of stepping-stone intrusion,” *expert systems with applications*, vol. 37, no. 2, pp. 1431–1437, 2010.
- [150] S.-H. S. Huang, H. Zhang, and M. Phay, “Detecting Stepping-stone intruders by identifying crossover packets in SSH connections,” in *2016 IEEE 30th International Conference on Advanced Information Networking and Applications (AINA)*. IEEE, 2016, pp. 1043–1050.
- [151] S.-H. S. Huang and Y.-W. Kuo, “Detecting chaff perturbation on stepping-stone connection,” in *2011 IEEE 17th International Conference on Parallel and Distributed Systems*. IEEE, 2011, pp. 660–667.
- [152] W. Ding, K. Le, and S.-H. S. Huang, “Detecting stepping-stones under the influence of packet jittering,” in *2013 9th International Conference on Information Assurance and Security (IAS)*. IEEE, 2013, pp. 31–36.
- [153] X. Wang and D. Reeves, “Robust correlation of encrypted attack traffic through stepping stones by flow watermarking,” *IEEE Transactions on Dependable and Secure Computing*, vol. 8, no. 3, pp. 434–449, 2010.
- [154] J. Yang and S.-H. S. Huang, “Mining TCP/IP packets to detect stepping-stone intrusion,” *computers & security*, vol. 26, no. 7-8, pp. 479–484, 2007.
- [155] X. Xu and X. Wang, “An adaptive network intrusion detection method based on PCA and support vector machines,” in *International Conference on Advanced Data Mining and Applications*. Springer, 2005, pp. 696–703.
- [156] R. R. Karthick, V. P. Hattiwale, and B. Ravindran, “Adaptive network intrusion detection system using a hybrid approach,” in *2012 Fourth International Conference on Communication Systems and Networks (COMSNETS 2012)*. IEEE, 2012, pp. 1–7.
- [157] J. Tetreault, J. Burstein, E. Kochmar, C. Leacock, and H. Yannakoudakis, “Proceedings of the Thirteenth Workshop on Innovative Use of NLP for Building

- Educational Applications,” in *Proceedings of the Thirteenth Workshop on Innovative Use of NLP for Building Educational Applications*, 2018.
- [158] A. Nourbakhsh and G. Bang, “A framework for anomaly detection using language modeling, and its applications to finance,” *arXiv preprint arXiv:1908.09156*, 2019.
- [159] Z. Chen, D. Chen, Z. Yuan, X. Cheng, and X. Zhang, “Learning Graph Structures with Transformer for Multivariate Time Series Anomaly Detection in IoT,” *arXiv preprint arXiv:2104.03466*, 2021.