

Examining traffic microstructures to improve model development

Henry Clausen
University of Edinburgh
Edinburgh, UK
henry.clausen@ed.ac.uk

David Aspinall
University of Edinburgh
Alan Turing Institute
Edinburgh, UK
david.aspinall@ed.ac.uk

Abstract—We demonstrate how machine-learning-based network intrusion detection models can be validated and developed by probing models using traffic with specifically controlled microstructures. We show our methodology by probing two published state-of-the-art models to find classification flaws and understand misbehaviour. These models fail for input traffic with particular characteristics such as retransmissions or overly dispersed flow interarrival times. After we make simple corresponding model corrections, detection rates already improve between 2–4%. We believe this shows promise for using tailored data with controllable and labelled characteristics to effectively improve model development in NID, a practice that helped model development significantly in several other areas of machine-learning.

Index Terms—Machine learning, traffic microstructures, network intrusion detection

I. INTRODUCTION

The model development process of machine-learning (ML) based network intrusion detection (NID) models usually ignores specific traffic characteristics and lacks the ability to extensively explore model failings. The main reason for this is likely the lack of precise datasets with specifically curated characteristics and corresponding information. In this paper, we demonstrate how the generation of traffic with controllable and labelled microstructures enables researchers to probe a model and its reaction to various traffic phenomena to much greater detail in order to understand and develop the model’s capabilities.

Machine-learning breakthroughs in other fields have often 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) [9] 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 understand these implicit structures

[5], which in turn leads to design improvements of general speech recognition models [8].

Prominent network intrusion detection methods as Kitsune [11] or DeepCorr [12] learn structures in the sizes, flags, or interarrival times of packets for decision-making. These **traffic microstructures** that can be observed when looking sequences of packets or connections reveal information about attack behaviour, but are also influenced by a number of other factors such as network congestion or the transmitted data. However, no effort has been made so far to monitor or control these factors to probe these models for specific microstructures, and researchers often just evaluate a variety of ML-models on sparsely labelled datasets in the hope of edging out competitors, without understanding model flaws and corresponding data structures through targeted probing.

In this work, we aim to demonstrate the usefulness of the control and information on traffic microstructures for model validation and development. We show the inspection of two state-of-the-art network intrusion detection models with specially generated traffic to identify model flaws, understand model behaviour better, and subsequently boost corresponding results. We hope to find new ways to improve model development in NID and increase the efficiency with which models can learn traffic microstructures.

A. Outline

The remainder of the paper is organized as follows. Section II discusses the necessity for probing to validate and understand ML-models, and our methodology of using traffic microstructure control for model probing. In Sections III and IV we demonstrate how to perform model probing and implement corresponding design improvements on two network intrusion detection models. Section V concludes the results and discusses limitations of our work and directions for future work.

II. MOTIVATION AND METHODOLOGY

A. Motivation

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**. Model

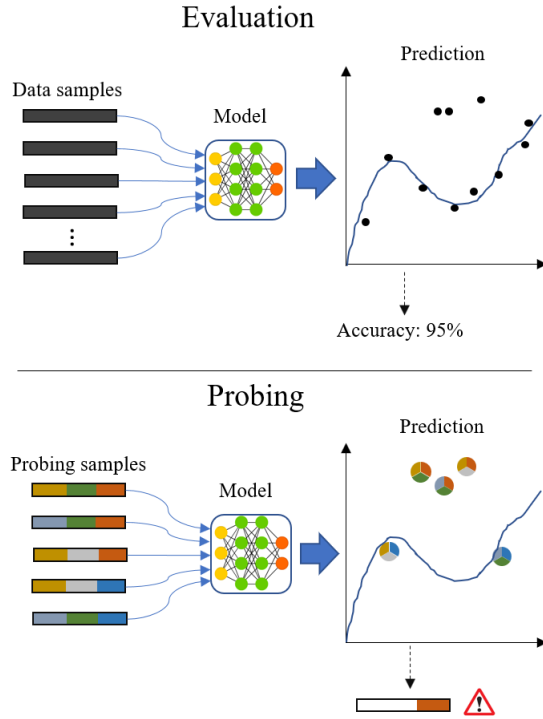


Fig. 1: Comparison between numerical model evaluation and model probing with specifically controlled data characteristics, indicated as colours.

validation is essential to understand how particular data structures are processed, and enables researchers to develop their models accordingly. The emergence of several perturbative model analysis tools that give insights into which structures are important in the decision-making of deep learning methods underline the importance of proper model validation. Existing tools such as the *What-If tool* [18] are however either only capable of computing perturbations for small data sequences, or rely on domain-knowledge such as *DeepLIFT* [15]. In the case of network traffic models, the lack of ground-truth labels for individual traffic traces prohibits us to associate patterns in a packet sequence with specific traffic shaping characteristics such as activities or congestion.

Recently, machine-learning methods have been trained on large datasets to classify raw traffic, often using simple malicious/benign labels. These models usually classify traffic traces as a sequence of packets within a connection (DeepCorr [12], Kitsune [11]), or as a short sequence of connections (Radford et al. [13]). Designs vary from models that make binary decisions by recognising characteristic patterns to models that identify anomalous traffic by learning generalised structures for all encountered benign sequences. We want to demonstrate how model validation can be performed for such models with the use of specifically generated traffic traces. We focus in particular on **traffic microstructures**, which we define as sequential structures in the packet or flow metadata stream. For this, we are generating traffic from different activities and settings in short time-intervals that are

suitable for models that classify **individual connections** or **short sequences** of connections, such as done by many traffic classifiers and low-volume intrusion detection systems. Models that depend on long-term or network wide structures, such as for the detection of botnet, DoS-attacks or fast-spreading worms do not fall within this scope.

B. Generating controllable traffic microstructures

We use DetGen [2], a tool that generates traffic with a fine-grained control over traffic shaping factors. In this work, we look traffic generated from scenarios that include regular HTTP communication, requests to an SQL-server, multi-host file-synchronisation, SQL-injection attacks, botnet traffic, as well as FTP-, SSH-, and SMTP-communication. DetGen offers control over the following traffic shaping factors:

- a) *Performed task/application*: The task and application drive the communication and influences characteristics such as direction and rate of transfer or the number of connections.
- b) *Application layer implementations*: Different implementations for TLS, HTTP, etc. can yield different channel prioritisation and can perform different handshakes.
- c) *Transferred data*: The amount and content of transferred data influences the overall packet number, rate, and size such as shown by Biernacki [1] for streaming services.
- d) *Caching/Repetition effects*: Tools like cookies, website caching, DNS caching, known hosts in SSH, etc. remove one or more information retrieval requests, which leads to altered packet sequences and less established connections [4].
- e) *Host level load*: Computational load (CPU, memory, I/O) on the host machine can affect the processing speed of incoming and outgoing traffic.
- f) *LAN and WAN congestion*: Low available bandwidth, long RTTs, or packet loss can have affect on TCP congestion control mechanisms and influence frame-sizes, IATs, and the overall temporal characteristic of the sequence.

In this work, we focus mainly on influence from the factors a), b), f) and partly c), as they have most influence on the behaviour of the models discussed in Section III and IV. The completeness of this list as well as the impact of the listed factors on overall model behaviour goes beyond the scope of this work and is better discussed by Clausen et al. [3].

Labels that describe the respective setting for each factor are attached to each traffic sample after generation, thus enabling us to provide ground-truth information about the precise generation setting of individual samples. Traffic is generated in a virtual network along with virtual software switches, Ethernet links and routers. The communication is mostly performed in a client-server setting, however some settings such as multi-host file-synchronisation involve more hosts.

C. Methodology

Before the probing, we have to identify which types of characteristics the model should be probed on, and generate the corresponding data. We then train the model mainly on the datasets that is used for the general evaluation, but also attach

probing data to make up at least 10% of the particular traffic type in the training set. This is to ensure that the model is able to see and learn the structures in the probing data, even though the overall type of traffic in the probing data should be similar to the evaluation data to provide a consistent model.

After training and general evaluation, the model probing is done by feeding the model data samples with the desired descriptive labels, monitoring the output or behaviours in dependence of these labels, and comparing them to the expected output or behaviour. Since each traffic sample contains multiple descriptive labels, it is possible to monitor the model response to multiple characteristics in parallel.

III. IMPROVED TRAFFIC SEPARATION FOR A CLASSIFIER WITH CONGESTION LEVEL INFORMATION

Our first example looks at how descriptive ground truth information on traffic characteristics can improve a traffic classification model through the analysis of data separation in dependence of different traffic features. For this, we use a *Long-Short-Term Memory* (LSTM) network, a deep learning design for sequential data, by Hwang et al. [7], 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 [16]. 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.

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 [14] (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. 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 [10].

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

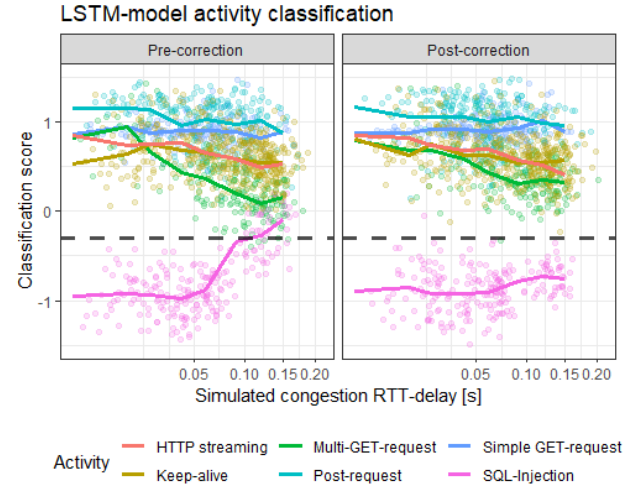


Fig. 2: Scores for the LSTM-traffic model before and after the model correction.

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).

Step 3: Fig. 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. 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

¹tuned for the geometric mean

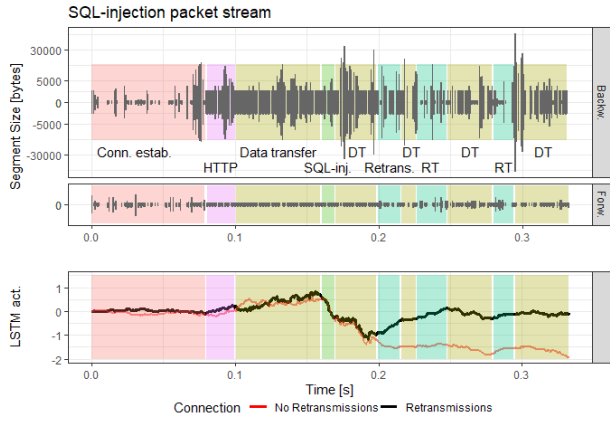


Fig. 3: LSTM-output activation in dependence of connection phases.

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. 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. 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%.

IV. 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* [11], a recent deep learning anomaly-detection model based on stacked autoencoders. *Kitsune*'s AUC-scores surpassed those of other state-

²scaled temporally to the same connection phases

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 I: 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.

of-the-art methods for a variety of attacks, including various types of Botnet traffic and *man-in-the-middle* attacks.

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 centers. 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 [14] 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.

A. 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 [6].

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 Bot-net 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 II-B being held constant within each setting except for small variations in the transmitted message or file. Table I summarises the traffic for each setting.

We verify if traffic samples within each group are projected to similar areas by measuring the average and maximum Mahalanobis-distance to quantify the overall dispersion of the samples. The results are displayed in Table I and depicted in Fig. 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 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.

B. 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. 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 HTTP-traffic with the label “Sudden

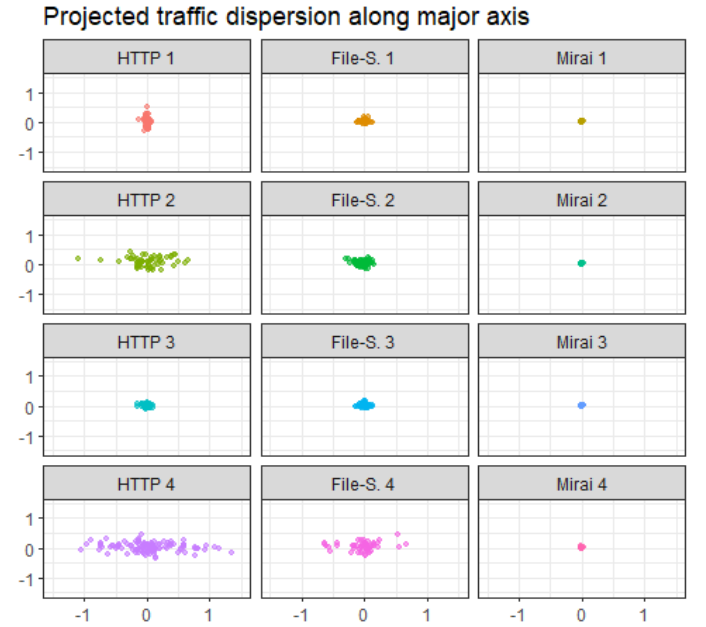


Fig. 4: Dispersion of projected traffic samples from each setting, plotted along the two most dispersed axes.

termination” are distributed across the cluster axis in a similar fashion, also depicted in Fig. 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%**.

V. CONCLUSIONS

In this paper, demonstrated the impact of traffic generation with extensive microstructure control as well as detailed corresponding documentation on researchers ability to evaluate and understand network intrusion detection models. We implemented and trained two state-of-the art detection models before extensively probing their behaviour and limitations when encountering different traffic types.

By using HTTP-traffic with congestion settings, we were quickly able to identify the inability of an LSTM-based classifier to handle traffic with significant retransmission rates, which enabled us to improve the model accordingly and increase detection performance by more than 2%. Similarly, the examination of projection consistency of a subspace-clustering method using traffic with artificially similar characteristics

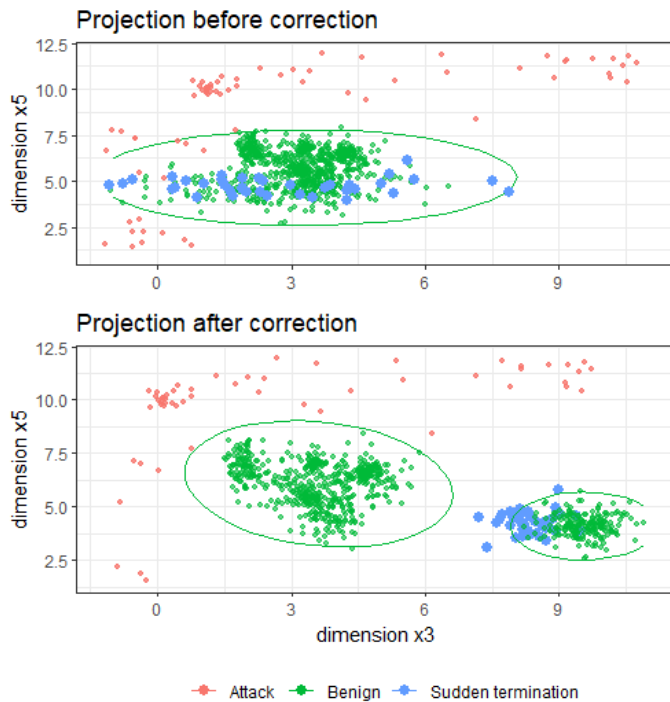


Fig. 5: Scores for the LSTM-traffic classification model in dependence of simulated network congestion, along with the classification threshold

revealed a high sensitivity to interarrival times, while cluster-coherence could be increased significantly by identifying half-open connections as the source of overly dispersed traffic projections.

We believe that in combination with strong NID-datasets, model probing with targeted traffic samples might hold the key to reduce false positives of detection models to an acceptable rate and help replicate detection rates in practical settings. This might furthermore help improve the transfer of learned structures from one to another dataset by identifying where the model is failing on the new data, and to compile a suitable dataset for fine-tuning training.

DetGen, the tool we used to control traffic microstructures with, is openly accessible on GitHub and is discussed in [2].

A. Difficulties and limitations

Controlling traffic shaping factors artificially can exaggerate some traffic characteristics in unrealistic ways and thus both affect the training phase of a model as well as tilt the actual detection performance of a model in either direction. Additionally, the artificial randomisation of traffic shaping factors can currently not generate the traffic diversity encountered in real-life traffic, something that is however significantly more pronounced in commonly used network intrusion datasets such as the CICIDS-17 dataset, where the majority of successful FTP-transfers download the Wikipedia page for ‘Encryption’.

Refining the ability of classifiers to identify traffic characteristics related to particular activities can lead to privacy

infringements or even discrimination against users or traffic types. It is therefore important to also investigate how applications can implement privacy-enhancing measures that conceal traffic characteristics, such as proposed by Wang et al. [17].

ACKNOWLEDGEMENT

The authors gratefully acknowledge funding support from BT Group PLC, UKRI/EPSC (grant EP/N510129/1) and the hosting University of Edinburgh.

REFERENCES

- [1] A. Biernacki. Analysis and modelling of traffic produced by adaptive http-based video. *Multimedia Tools and Applications*, 76(10):12347–12368, 2017.
- [2] H. Clausen, R. Flood, and D. Aspinall. Traffic generation using containerization for machine learning. In *Dynamic and Novel Advances in Machine Learning and Intelligent Cyber Security (DYNAMICS) Workshop*, 2019.
- [3] H. Clausen, R. Flood, and D. Aspinall. Controlling network traffic microstructures for machine-learning model probing. 2021. Manuscript submitted for publication.
- [4] 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*, pages 310–315. IEEE, 2012.
- [5] 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)*, pages 7355–7359. IEEE, 2019.
- [6] 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)*, pages 1133–1141. IEEE, 2017.
- [7] 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*, 9(16):3414, 2019.
- [8] 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)*, pages 6414–6418. IEEE, 2020.
- [9] 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*, 13(5):e0196391, 2018.
- [10] P. Mell, V. Hu, R. Lippmann, J. Haines, and M. Zissman. An overview of issues in testing intrusion detection systems, 2003.
- [11] 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.
- [12] 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*, pages 1962–1976, 2018.
- [13] 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.
- [14] I. Sharafaldin, A. H. Lashkari, and A. A. Ghorbani. Toward generating a new intrusion detection dataset and intrusion traffic characterization. In *ICISSP*, pages 108–116, 2018.
- [15] A. Shrikumar, P. Greenside, and A. Kundaje. Learning important features through propagating activation differences. In *International Conference on Machine Learning*, pages 3145–3153. PMLR, 2017.
- [16] 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*, 154:102538, 2020.
- [17] T. Wang and I. Goldberg. Walkie-talkie: An efficient defense against passive website fingerprinting attacks. In *26th {USENIX} Security Symposium ({USENIX} Security 17)*, pages 1375–1390, 2017.
- [18] 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*, 26(1):56–65, 2019.