

# TANTRA: Timing-Based Adversarial Network Traffic Reshaping Attack

Yam Sharon

*Ben Gurion University of the Negev*

Yang Liu

*Nanyang Technological University*

David Berend

*Nanyang Technological University*

Asaf Shabtai

*Ben Gurion University of the Negev*

Yuval Elovici

*Ben Gurion University of the Negev*

## Abstract

Network intrusion attacks are a known threat. To detect such attacks, network intrusion detection systems (NIDSs) have been developed and deployed. These systems apply machine learning models to high-dimensional vectors of features extracted from network traffic to detect intrusions. Advances in NIDSs have made it challenging for attackers, who must execute attacks without being detected by these systems. Prior research on bypassing NIDSs has mainly focused on perturbing the features extracted from the attack traffic to fool the detection system, however, this may jeopardize the attack's functionality. In this work, we present TANTRA, a novel end-to-end **T**iming-based **A**dversarial **N**etwork **T**raffic **R**eshaping **A**ttack that can bypass a variety of NIDSs. Our evasion attack utilizes a long short-term memory (LSTM) deep neural network (DNN) which is trained to learn the time differences between the target network's benign packets. The trained LSTM is used to set the time differences between the malicious traffic packets (attack), without changing their content, such that they will "behave" like benign network traffic and will not be detected as an intrusion. We evaluate TANTRA on eight common intrusion attacks and three state-of-the-art NIDS systems, achieving an average success rate of 99.99% in network intrusion detection system evasion. We also propose a novel mitigation technique to address this new evasion attack.

## 1 Introduction

The networks of governments, enterprises, and private households are the target of network intrusion attacks performed by attackers with an interest in financial gain, disruption, or gathering intelligence. According to recent trends, the most prominent reason for targeted network intrusion attacks is acquiring intelligence. Such attacks increased tenfold from 2017 to 2018 [34]. In order to launch an attack on an organizational network (for example, a man-in-the-middle (MITM) attack), the attacker must first bypass the targeted network's defense mechanisms.

The defense mechanisms aimed at thwarting such attacks are referred to as network intrusion detection systems (NIDSs); to address the increasing complexity and sophistication of attacks, these systems have become more advanced in recent years. At the core of these defense advancements lies the extraction of high-dimensional vectors of features from network traffic, which are then utilized by traditional machine learning or deep learning models trained to distinguish between benign and malicious traffic.

Advances in NIDSs have presented attackers with a challenge which they have shown limited ability in overcoming [7, 29, 36]. The majority of evasion attack approaches are based on perturbing the attack's traffic features by assuming a white-box setting in which the attacker has knowledge about the system's feature extractor. This knowledge is then used to perturb malicious traffic by performing an adversarial attack aimed at bypassing the NIDS. While shown to be effective at bypassing the NIDS, these attacks lose some of their original functionality after the adversarial perturbation is applied. Thus, the attacks' impact on the target network remains limited, as the attacks may not be viable in an end-to-end setting.

Given the limitations of feature-based attacks, new ways of evading NIDSs are likely to be explored by attackers. One promising direction aims at enabling malicious intrusion traffic to mimic benign network behavior. In one of the first works aimed at mimicking benign traffic, Han et al. reshaped the timing of malicious network traffic and added dummy packets to avoid detection by the NIDS [16]. This approach requires preparing all of the attack traffic on the attacker's side and was not evaluated in an end-to-end scenario. Instead, the performance was measured against a NIDS, without setting up an active network connection. Given the need to prepare all of the attack traffic on the attacker's side in advance, the approach cannot be adjusted according to frequent target network responses or delays. This makes it less applicable for real-world settings.

In this paper, we present a new **T**iming-based **A**dversarial **N**etwork **T**raffic **R**eshaping **A**ttack (TANTRA) that mimics

benign traffic to evade detection. TANTRA’s main advantages are that it maintains the attack packets’ content and can be adjusted according to frequent target network responses in real time. Furthermore, no knowledge on the NIDS is required, resulting in a fully black-box attack. This makes TANTRA the first successful real-world evasion attack against NIDSs.

The proposed evasion attack is based on training a long short-term memory (LSTM) to generate the time differences between benign network packets (collected from the target network in advance). Then, the trained model is used to reshape the malicious traffic’s inter-packet delay. This way, the modified malicious traffic is able to mimic the behavior of benign traffic. We evaluate TANTRA on eight common network attacks against three state-of-the-art (SOTA) NIDSs that are based on traditional machine learning (ML) and deep learning (DL) techniques. Our novel evasion attack successfully evaded all three defenses, achieving a 99.99% success rate (on average), outperforming evasion attacks proposed in prior research under similar settings by up to 28.8% (as shown in Section 6).

In summary, our contributions are as follows:

- We present TANTRA, a novel network evasion attack that can evade SOTA defenses with a 99.99% success rate. The proposed attack outperforms methods/attacks proposed in related work and performed under similar conditions by up to 28.8% more effectively.
- TANTRA can evade NIDSs without changing the attack packets’ content, eliminating the risk of losing the attack’s functionality, which has been a prominent drawback in related work.
- TANTRA is the first end-to-end evasion attack against NIDSs that operates in a black-box setting, as successfully demonstrated against three SOTA NIDSs in this paper.
- To mitigate the threat posed by this new evasion attack, we propose a novel defense technique. A preliminary version of this mitigation was implemented, evaluated, and shown to be promising.

The remainder of the paper is structured as follows. In Section 2, the background is presented; this is followed by the threat model, which is described in Section 3. The attack methodology is discussed in Section 4, followed by a discussion of the weaknesses and mitigation of NIDSs in Section 5. We present the evaluation in Section 6 and provide a detailed comparison to related work in Section 7. Our conclusions are presented in Section 8.

## 2 Background

### 2.1 Network Intrusion Detection Systems

Network intrusion attacks refer to any unauthorized activity on a target network. To detect and prevent such attacks, NIDSs are employed. These systems are based on two main approaches: signature-based and anomaly-based approaches [4]. The signature-based approach is a lightweight defense where the defender manually defines a set of rules that characterize malicious traffic. This approach is easy to implement but is also easy to bypass, as the attacker can determine the rules and reshape the malicious traffic in light of the rules defined by the defender. It becomes more difficult for the attacker when the anomaly-based approach is utilized. When combined with deep learning techniques, this defense has been shown to be very effective [7, 18]. Together with the high-dimensional features extracted from the traffic, the NIDS receives a fine-grained perspective on the traffic characteristics; therefore, the NIDS can more effectively detect incoming attacks as anomalous. Hence, reshaping the malicious traffic so that it will not be detected as anomalous is challenging.

### 2.2 Traffic Reshaping

Given the effectiveness of recent NIDSs, attackers will be forced to explore new ways of performing network intrusion attacks. To be successful, the network intrusion attacks require adjustments to prevent detection by any NIDS. Modifying the network traffic packets of the attack is referred to as traffic reshaping; there are several approaches used for reshaping. One approach is to reshape the features extracted from the network traffic (feature-space reshaping). Another emerging direction is to reshape the timing of malicious network packets (timing-based reshaping).

For feature-space reshaping, attackers require knowledge of the target network feature extractor and NIDS, which is difficult to obtain. In this case, perturbations are applied to the extracted features in order to bypass the NIDS. These perturbations are specifically crafted for the targeted NIDS. The perturbations are referred to as adversarial examples, which, in the context of network intrusion, result in adversarial attacks. Adversarial attacks are inspired by the computer vision (CV) domain, where Szegedy et al. [35] first created a specially crafted pixel mask invisible to the human eye, causing the neural network to misclassify the input image. Later, Goodfellow et al. [14], Carlini and Wagner [6], and others [8, 26] advanced the field by minimizing the adjustments needed for a white-box setting. For intrusion attacks, such a setting requires extensive knowledge of the preprocessing techniques and behavior of the targeted NIDS. Previous studies have demonstrated successful intrusions when assuming a white-box setting [1, 10, 20, 30, 37–39]; however when a white-box setting is required, applicability in the real world

Table 1: Comparison to prior research.

Assumptions	Related work	Ours
Ongoing connection to target network	✓ [10, 20, 30, 37–39]	✓
Knowledge on feature extractor	✓ [10, 20, 30, 37–39]	.
Ability to change extracted features	✓ [10, 30, 39]	.
Access to NIDS	✓ [10, 30, 39]	.

is minimal. For greater real-world applicability, a black-box setting (commonly referred to as an end-to-end scenario) is required. For this scenario, knowledge and access to the targeted network’s feature extractor and feature space are not available. Furthermore, for the end-to-end scenario, the attacker must be capable of adjusting the evasion attack according to frequent responses of the target network. Finally, the attack’s functionality must be maintained.

Timing-based reshaping is emerging as an approach that does not require any knowledge of the NIDS. Its only requirement is the need to obtain benign traffic from the target network. This makes it a suitable candidate for a black-box setting. In [16] the authors used a generative adversarial network (GAN) to combine dummy packets with the malicious network packets to change the timing of the attack’s traffic. While showing promise for a scenario in which there is only limited knowledge on the target network, an end-to-end scenario was not able of being performed. To accomplish this, one remaining challenge that needs to be addressed lies in its methodology, where target network responses cannot be integrated. So far, the attack traffic is prepared (in its entirety) in advance which prevents response integration. Furthermore, the additional dummy packets pose a risk to the attack’s functionality. We aim to address such drawbacks by demonstrating a full end-to-end attack scenario without risking the attack’s functionality.

### 3 Threat Model

In this study, we assume that (1) the attacker treats the NIDS as a black-box, (2) the attacker has access to the target network, enabling him/her to observe the benign traffic (to learn its patterns) and send the reshaped attack traffic to the network, (3) the attacker adjusts the malicious traffic in real time, and (4) the attacker maintains the malicious packets’ content. These assumptions represent a realistic threat model and to the best of our knowledge, were not addressed in prior work.

During the setup phase, the attacker learns the benign timing behavior of the target’s network traffic. Then, the malicious network traffic is reshaped in real time, in accordance with the learned benign timing behavior, so that the target NIDS will not issue an alert. One of the main strengths of the proposed threat model is the limited assumptions necessary

for its execution, as presented in Table 1, which compares our assumptions to those of prior studies [10, 20, 30, 37–39]. In our approach, the attacker does not need access to the targeted NIDS or knowledge about its behavior (e.g., extracted features, NIDS type, architecture). The only requirement for the attacker is an ongoing connection to the target network, a requirement for any network intrusion attack.

## 4 TANTRA: Novel Evasion Attack

### 4.1 Main Challenges

An effective evasion attack against a NIDS was demonstrated in prior research, however, there were difficulties in retaining the attack’s functionality, and a white-box setting was required. The methodology used in our study reflects our aim of crafting an attack capable of both maintaining the attack packets’ content and evading the NIDS, with limited available knowledge. To achieve this we have to overcome two main challenges. The first challenge is switching the reshaping context from the commonly used feature space to the less-explored problem space. The second challenge is utilizing the limited network characteristics in the problem space for an effective evasion attack in a black-box setting.

When analyzing other approaches proposed for evading NIDSs [1, 10, 30, 37–39], we see that retaining the functionality of the attack packets’ content does not seem possible. In fact, Pierazzi et al. [1] found that modifying the feature space after feature extraction may lead to a loss in functionality. Analyzing their results, we found that extracted features are based on statistical functions, such as the mean and STD, which are difficult (if not impossible) to reverse once a perturbation has been added.

Furthermore, the extracted features are related to each other and calculated from the same network packets’ attributes. Changing one feature without changing the other related high-dimensional features will make it even more problematic to reverse the traffic to its original form. Therefore, we conclude that the feature space does not seem like a promising avenue to pursue for crafting an effective attack. Instead, we turn to the problem space. Here, reshaping can be done before the features are extracted from the traffic by the NIDS, and therefore it provides a foundation for retaining the functionality of the attack in a black-box setting.

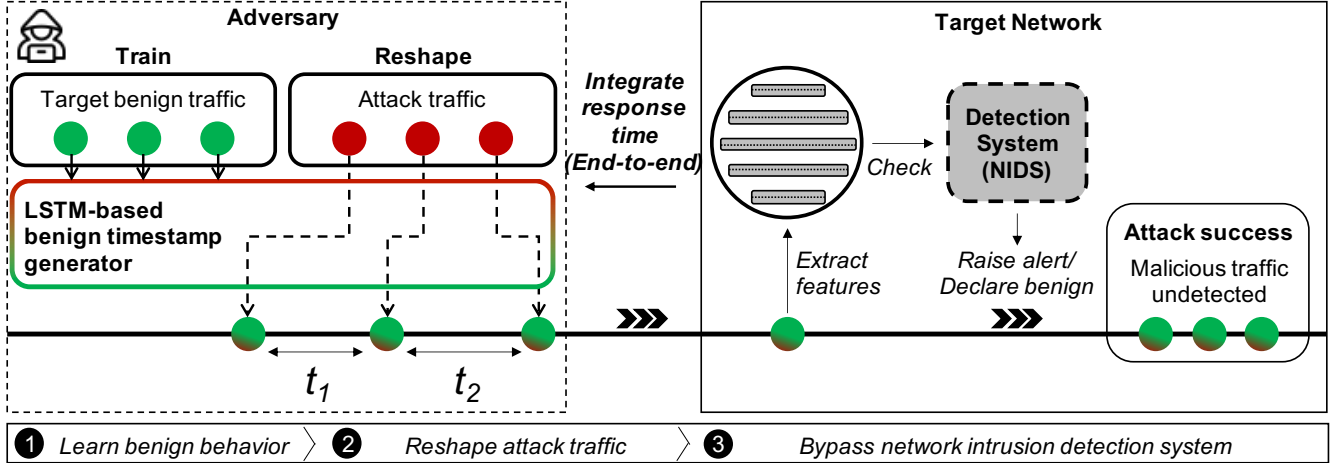


Figure 1: Overview of our evasion attack.

However, the problem space offers limited options for reshaping the network traffic characteristics. For example, most of the time, the attacker is unable to control or modify the network connection characteristics, such as the destination’s MAC and IP addresses. Two of the limited modifiable options are the packet size and timestamp features; the attacker can split or pad the packets’ content, and he/she can send the packets with different timestamps. In this work, we show that successful evasion can be achieved by reshaping just the packet timestamps, without influencing the packet’s size.

## 4.2 Overview

Our evasion attack consists of three steps, which are illustrated in Figure 1. In step ❶, an LSTM is trained to learn benign network traffic behavior by predicting the time differences between benign network packets. The LSTM is trained for the specific network connection on which the attack will be executed. In step ❷, the trained LSTM is applied on an intrusion attack’s malicious traffic to reshape the interpacket delay (adjusting the timestamps), so it is similar to that of benign traffic. In step ❸, the reshaped malicious traffic is sent to the target network where it aims to bypass the NIDS which is treated as a black-box.

## 4.3 Preparation

To reshape malicious traffic so it behaves like benign traffic, we utilize a many-to-one LSTM model. The LSTM’s input is a short history of previous packets’ timestamps, along with the size of the previous and current packets. To determine the optimal amount of history of network packets (window size), we perform an empirical assessment (see Section 6), starting with a window size (WS) of three and increasing it to 150. The LSTM’s output is the timestamp for the current network

packet. The architecture (Figure 2) consists of an LSTM cell for each input packet  $p$ , with a hidden layer  $h$  and cell output  $o$ . The cell output is passed on to a  $32 \times 8$  dense layer with ReLU activation. Finally, an  $8 \times 1$  dense layer with sigmoid activation generates the timestamp  $t$  (this architecture is based on previous work that proposed a many-to-one LSTM [5, 31]).

$$f(\text{history}, p_i) = t_i \quad (1)$$

During the training phase, the LSTM receives the benign traffic collected from the targeted network. In each training iteration, the LSTM receives  $n$  packets as input. Afterward,  $n$  timestamp predictions are compared with the known interpacket delays to compute the mean squared error (MSE) loss. The loss is then backpropagated through the LSTM. After the training phase, the LSTM is assumed to be capable of predicting interpacket delays between the current and previous packets using timestamps. One advantage of the LSTM architecture is its sequential behavior, where the traffic is processed packet by packet. This enables an end-to-end scenario.

## 4.4 Execution

The trained LSTM can be utilized to predict the timestamps of the attack, packet by packet. This enables an end-to-end scenario in which target network delays and responses are integrated during the attack reshaping process.

Since TANTRA is the first evasion attack capable of performing an end-to-end scenario, comparing the performance in this setup to related work may be inaccurate. Hence, we first analyze its performance directly against the NIDS, without an active target network connection, following the setup of prior research [10, 20, 30, 37–39]. In this scenario, malicious network traffic stemming from, e.g., an MITM attack, is

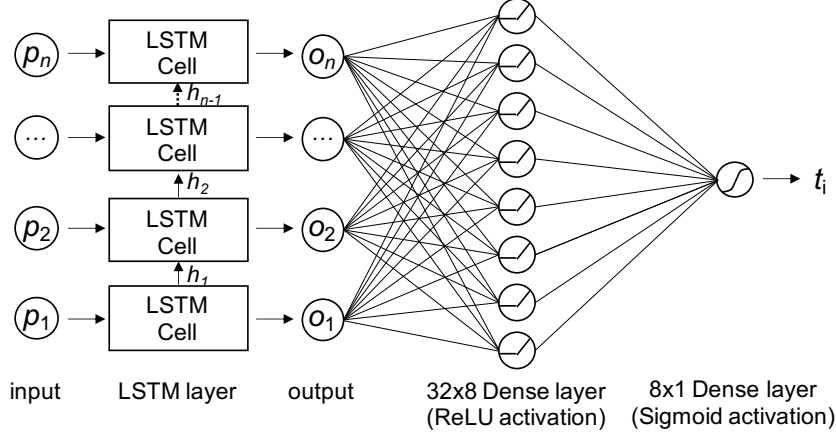


Figure 2: Overview of the DNN model architecture.

reshaped by adjusting the timestamps and evaluated directly against the NIDSs.

For the end-to-end scenario, the timestamps are generated step by step with an active target network connection. Here, the MITM attack is not directly evaluated on the NIDS; instead, it is sent in multiple steps to the target network. When the packets reach the target network, they are processed step-wise by the NIDS. Each step that follows integrates the target network delays and responses. For target network response integration, TANTRA is executed on a proxy that controls the timing of the outgoing traffic, packet by packet. Therefore, packets can be adjusted in real time. This represents a much more challenging setup, but it enables a real-world applicable execution, as TANTRA adjusts its approach according to frequent target network’s response.

#### 4.5 Novel Defense Directions

For mitigation, we look to the principles of CV-based DeepFake detection. DeepFake primarily relates to content, such as videos, generated by deep learning networks [24]. DeepFake algorithms can be used for malicious activities that may cause psychological, political, monetary, and physical harm. Given such implications, defenses have been developed using machine learning models to differentiate between real and fake videos (DeepFakes) [13, 15].

A common approach used is to adjust the training procedure of existing anomaly detectors so that they are aware of the existence of DeepFakes. Similarly, we propose training three different ML classifiers (NIDS), namely logistic regression, Gaussian naive Bayes, and random forest, to learn to detect reshaped traffic. We hypothesize that increasing NIDS’ awareness of reshaped characteristics will enable them to distinguish between benign traffic and reshaped malicious traffic. More information on NIDS weakness is described in the Section 5.

### 5 Weaknesses of NIDS

We observe that existing NIDSs aim to detect malicious traffic by searching for characteristics which are considered abnormal for benign traffic. This approach is used by all three state-of-the-art NIDSs examined in this study. Therefore, each NIDS is trained on benign traffic and searches for anomalies in network traffic to identify attacks.

For higher granularity in their prediction, all three of the systems perform their detection based on the feature space. Common practices involve using the AfterImage [23] feature extractor, which achieves SOTA results, represents the industry standard [7, 16, 32], and is used in this work. Using AfterImage, nine basic problem-space attributes from each packet (MAC, IP, and protocol attributes for each source and destination, together with the IP type, packet size, and timestamp) are used to extract features. Modifying one attribute affects all other attributes in the feature space. Therefore, changing the timestamp attribute affects the other attributes in the feature space. The extracted features then serve as input to the NIDS. Therefore, TANTRA aims to exploit the fact that NIDSs are trained only on benign traffic and reshapes malicious traffic accordingly. We assess this weakness for each NIDS individually below.

**Autoencoder (AE).** AEs are neural networks that aim to reconstruct their input; in doing so, AEs attempt to characterize specific behavior of and patterns in the traffic, enabling them able to distinguish differences in high-dimensional space when abnormal input is introduced. The AE represents a function  $f$  which aims to deconstruct and reconstruct an input. For any given input  $x$  from the training set,  $f(x)$  aims at minimizing the reconstruction error RMSE. In the context of intrusion detection, the AE aims for the benign network packet,  $b$ , and its output  $f(b)$ ,  $f(b) - b = RMSE$  to be small. Similarly, it aims for a malicious network packet,  $m$ , and its output  $f(m)$ ,  $f(b) - b = RMSE$ , to be large. In our case, we use the com-

monly used root mean square error (RMSE) as reconstruction error.

$$RMSE(\vec{X}, \vec{Y}) = \sqrt{\frac{\sum_{i=1}^n (x_i - y_i)^2}{n}} \quad (2)$$

Our novel reshaping technique exploits this type of NIDS by modifying the malicious network traffic so it follows benign behavior. Since it reshapes the timing characteristics of the attacker’s malicious traffic, the attack will be perceived similarly to benign traffic by the AE. Thus, we hypothesize that the reconstruction error (RMSE) between  $m$  and  $f(m)$  will be both small and below a defined detection threshold.

**KitNET.** Proposed by Mirsky et al. [23], KitNET is an unsupervised neural network composed of an ensemble of AEs. Using an ensemble of AEs provides the ability to use simpler AE architectures and lowers the computational complexity. Therefore, like an AE, KitNET tries to characterize benign network traffic, determining  $f(b) - b = RMSE$  that knows how to reconstruct benign network traffic  $b$ . We expect KitNET’s reconstruction error of reshaped malicious traffic to be small, like that of AEs, and therefore likely to evade detection.

**Isolation Forest.** The isolation forest (IF) [21] algorithm is an anomaly detection algorithm commonly used in related studies [11, 22, 28]. IF creates clusters of benign data; when encountering malicious input, the IF algorithm should isolate such inputs, based on their features, outside the clusters. IF defines its anomaly score for a given example  $x$  as:

$$s(x, n) = 2^{-\frac{E(h(x))}{c(n)}} \quad (3)$$

where  $E(h(x))$  is the average path length of observation,  $x$ , to a cluster (referred to as a tree);  $c(n)$  is the average path length of unsuccessful searches in the trees; and,  $n$ , is the number of external nodes. An anomaly score close to one indicates an anomaly. Scores much smaller than 0.5 indicate normal observations, as the path’s length to the current example,  $x$ , is smaller or equal to the average path length. The theory behind the IF model includes the principle that anomalous examples behave differently than benign ones; thus, it will be easier to isolate anomalous examples than to detect examples which are typical for benign cases. After reshaping the malicious traffic so it behaves similarly to benign traffic, we hypothesize that  $E(h(x))$  will be close to the path length of benign examples,  $c(n)$ . Therefore, the reshaped malicious traffic should be similar to benign traffic and thus will have a similar average path length and evade detection.

## 6 Evaluation

In this section, we evaluate TANTRA’s effectiveness and mitigation. We first examine its baseline performance against

state-of-the-art NIDSs following the setup used in prior research (Section 5). Then, we optimize the attack by identifying the optimal LSTM window size and assess the attack’s effectiveness in an end-to-end scenario with an active network connection. Finally, we evaluate the proposed defense technique as a mitigation against our novel evasion attack.

### 6.1 Experimental Setup

#### 6.1.1 Datasets

We use the Kitsune dataset [23] and intrusion detection evaluation dataset CIC-IDS2017 [33]. Both benchmarks have been used in related studies [10, 16]. The datasets contain network traffic recordings of eight real network attacks, along with the benign and malicious labels for the network traffic packets. The network traffic recordings are provided as PCAP files, which are retrieved in individual network setups. The overall length of the recordings ranges from 28.15 to 118.9 minutes, with an average of 59.35 minutes, and an average of 2,540,013 packets. Table 2 presents (1) the total recording time, (2) the portion used to train the NIDS (benign traffic only), (3) the portion used to train TANTRA’s LSTM (benign traffic only), and (4) the portion used to evaluate TANTRA (malicious traffic only).

#### 6.1.2 NIDS

TANTRA is evaluated on three NIDSs that cover both traditional machine learning and deep learning approaches. After the traffic is reshaped by TANTRA in the problem space, it is extracted to the feature space where it serves as input to the NIDS. As done in other studies, we use the AfterImage feature extractor which converts the PCAP files into features that can serve as NIDS input.

**Autoencoder.** The AE NIDS consists of two parts - an encoder and a decoder. We use a convolutional neural network with a two-layer encoder (100x64 neurons), latent space of 32 neurons, and two-layer decoder architecture (64x100 neurons). For training, the Adam optimizer is utilized with a learning rate of 0.001, and the mean squared error (MSE) is used as the loss criterion. The Keras framework is employed [9].

**KitNET.** We follow KitNET NIDS’ open-source code with the default maximum autoencoder size of 10 in the ensemble layer, a learning rate of 0.1, and a hidden ratio of 0.75. We set the mapping and training phase parameters according to the work presenting Kitsune. RMSE is used as the loss criterion.

**Isolation Forest.** The Isolation Forest NIDS has 100 base estimators in the ensemble, with the threshold set to auto contamination, and the random state equal to zero and default anomaly score function [27].

Table 2: The datasets and utilization for experiments (time in minutes).

Dataset	Attack Type	Total Recording	Benign NIDS Training	Benign LSTM Training	Malicious Testing
Kitsune Dataset [23]	Active Wiretap	21.93	8.27	8.12	11.46
	MITM	20.17	8.05	7.16	6.97
	Fuzzing	28.95	12.96	10.91	14.58
	Mirai	118.95	33.68	31.61	75.82
	SSDP Flood	40.74	14.40	14.37	4.04
	SSL Renegotiation	38.73	14.69	12.38	22.05
CIC IDS 2017 [33]	Brute-Force	61.1	14.43	0.23	4.15
	SQL Injection	57.25	14.43	11.28	2.8

### 6.1.3 Evaluation Metrics

The detection rate (DR) of the NIDS for detecting malicious traffic is used to evaluate TANTRA’s performance. The changes in the DR as a result of our evasion attack are used to quantify TANTRA’s effectiveness; this is also referred to as recall. The DR metric has typically been used to showcase either a large DR when proposing a new NIDS [12, 19, 23] or a low DR when proposing a new attack technique [20, 37, 39].

## 6.2 TANTRA Baseline Evaluation

In this section, we provide a baseline evaluation, following prior work in which all of the traffic is prepared on the attacker’s side and evaluated against the NIDS directly. This evaluation setup disregards any potential target network behavior, such as network delays and responses, which may influence successful evasion. Although past studies have been unable to demonstrate an effective evasion attack in an end-to-end scenario, we perform the first step of the evaluation using their setup for comparison and method optimization purposes.

We train an LSTM with the window size set at three, which also utilizes an Adam optimizer with 0.001 learning rate, and the mean squared error (MSE) loss function. For each attack, an LSTM is trained on benign traffic that matches the traffic on the targeted on which the attack is performed. On average, the benign traffic (training set) for each attack consists of 296,501 packets from the network recordings, with a minimum of 73,875 packets. We evaluate the DR of the three NIDSs before and after reshaping the attack network traffic.

The results presented in Table 3 show that after applying our evasion attack to each malicious dataset, the network traffic in each dataset is able to successfully evade the NIDS (lowering the DR). After reshaping the attack traffic using our evasion attack, the DR decreases on average from 61.19% to 2.06%.

Without using TANTRA, the AE NIDS is able to detect most network attacks with an average DR of 61.11%. However, after reshaping the attack, TANTRA almost completely conceals the attack, reducing the average DR on the AE to

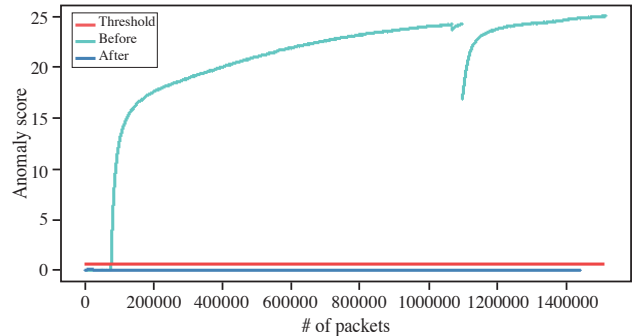


Figure 3: AE NIDS anomaly score for Mirai attack.

1.55%. While the DR of most attacks is reduced to below 1%, the SSDP flood attack is an exception. In this case, the initial detection rate of 71.94% is reduced to 10.22%, which is ten times higher than the average DR after reshaping for the other attacks. One possible explanation for this behavior is that the SSDP attack is an attack that contains several phases, with responses between some of the phases. Since all of the attack network traffic is generated before being sent to the target network, this evaluation setting prohibits integrating the network response time in each phase. Thereby, this illustrates the drawback of the evaluation setting used in past work and the relevance of evaluating an end-to-end scenario.

Figure 3 presents an example of the behavior of the AE when classifying malicious traffic. In the example, we present the anomaly score for the Mirai attack, before and after using TANTRA. The x-axis represents the number of processed network packets in a sequence, and the y-axis represents the anomaly score. The anomaly score before TANTRA (turquoise) is far above the detection threshold (red line). After TANTRA, the anomaly score for the Mirai attack is below the threshold for all network traffic. Therefore it can be seen that the AE is unable to detect the malicious behavior, which results in a DR of 0%.

The greatest impact of the evasion attack is observed when the KitNET NIDS is used. Before reshaping, KitNET performs best, with an average DR of 70.61%. After reshaping,

Table 3: NIDS detection rate before and after TANTRA (in %).

Attack Type	AE		KitNET		IF		KitNET (After)	
	Before	After↓	Before	After↓	Before	After↓	Related Work [16]↓	TANTRA↓
Active Wiretap	98.03	0.00	99.02	0.00	68.73	7.29	N/A	N/A
MitM	23.79	0.00	78.20	1.37	1.47	0.00	N/A	N/A
Fuzzing	67.53	0.97	92.37	0.00	49.28	0.83	1.31	0.00
Mirai	88.94	0.00	100.00	0.00	0.83	0.42	0.58	0.00
SSDP Flood	71.94	10.22	99.97	0.00	99.94	0.01	21.47	0.00
SSL Renegotiation	89.34	1.19	39.19	0.00	6.80	0.00	N/A	N/A
Brute-Force	25.50	0.00	22.79	0.00	87.94	9.34	28.18	0.00
SQL Injection	23.81	0.00	33.33	0.00	100.00	17.91	N/A	N/A
<b>Average</b>	61.11	1.55	70.61	0.17	51.87	4.48	12.89	0.00

the average DR decreases to 0.17%, the lowest DR among the three NIDSs. Another interesting finding is that TANTRA achieves a perfect score for seven of the eight network attacks examined. One possible explanation for the significant results when using KitNET is that the autoencoder ensemble is likely to be the best at characterizing benign behavior. As TANTRA aims to exploit benign behavior, it reshapes the malicious traffic so it **behaves** like benign traffic, thereby evading detection, aligned with our theoretical assessment (see Section 5).

Prior research has also used KitNET NIDS for evaluation [16], which serves as baseline comparison. Here, the performance is evaluated on four attacks, as shown in Table 3. The results of the study show that the attacks presented in prior research are able to decrease detection by the KitNET NIDS by an average of 12.89% with DRs ranging from 0.58% to 28.18%. For the evaluated attacks, TANTRA decreases detection to 0.00% for all attacks. Other prior research used white-box settings which are difficult to compare directly to TANTRA. Therefore we present additional qualitative comparison in Section 7.

Figure 4 provides a visual example for KitNET. This figure presents the anomaly score for the SSDP flood attack. Again, the turquoise lines, which represent the malicious traffic before reshaping, are far above the detection threshold. However, after reshaping, the NIDS fails to detect the traffic as malicious, as represented by the dark blue lines below the red threshold.

The lowest DR before TANTRA is applied is obtained by the IF NIDS. Before reshaping, the IF has an average DR of 51.87%, obtaining a DR below 10% with three of the eight attacks. However, after reshaping, the IF has the highest DR (an average DR of 4.48%) of the three intrusion detection systems. This DR stems from the relatively higher detection rate this NIDS obtained for three attacks after TANTRA was used, namely SQL injection with a DR of 17.91%, brute-force with 9.34%, and active wiretap with 7.29%. The DR for the other five attacks is below 1.0%.

Figure 5 presents the anomaly scores of the IF for the MITM attack, before and after TANTRA. The results support our claim that many malicious packets are different from

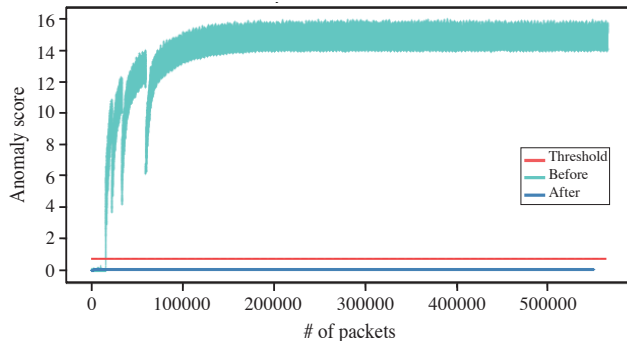


Figure 4: KitNET NIDS anomaly scores for SSDP flood attack.

benign examples before reshaping. After reshaping, most malicious packets are labeled benign by the IF model.

### 6.2.1 LSTM Optimization

While improving the attacks’ ability to evade detection, for most attacks the use of TANTRA did not result in a DR of 0%. For example using TANTRA for the active wiretap attack against the IF NIDS results in a DR of 7.29%. When using TANTRA for the SSDP flood attack against the AE, the DR is 10.22%. Therefore, in this section, we explore how TANTRA can be optimized, to ensure that it is capable of fully evading any type of state-of-the-art NIDS, regardless of the type of attack being reshaped.

Given the problem space’s limited attributes, we turned to the LSTM hyperparameters for optimization. Here, we focus on the window size, which controls how much information about past predictions influences the current one. For example, with a window size of three, only the last three prediction outcomes are integrated into the prediction, while with a window size of 10, the last 10 predictions are utilized, and so on.

There is no perfect solution for the optimal window size, as this varies depending on the attack and input size. Olah et al. [25] demonstrated this issue in the natural language processing (NLP) field with a translation task. Given the



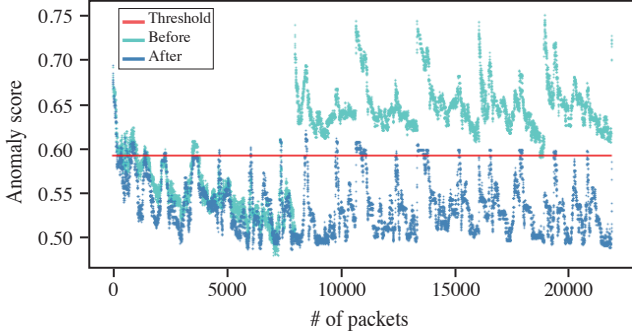


Figure 5: IF NIDS anomaly score for MITM attack.

sentence, "I grew up in France. I speak fluent French," our goal is to predict the last word. If we only take the previous three words into account, "I speak fluent \_\_\_\_\_," our model will most likely not predict the right word, whereas if we use the previous four words or more, the probability of the model to be correct is much higher, as the word "France" is now part of the history.

Factors to consider regarding the appropriate window size are the overall length of the input (analogous to the length of a sentence) and the context dependencies between the input elements (analogous to the context dependencies between words in a sentence). For TANTRA's LSTM, the input elements are network packet properties. The length of the network traffic varies depending on the attack. In the datasets, an active wiretap attack consists of roughly 500,000 packets, while the brute-force attack has only 1,507 packets. Therefore, a different window size may be needed depending on the attack.

The results presented in Table 4 show, that increasing the window size improves the attack's evasion capabilities, decreasing the average DR to 2.98% for WS=3, 0.88% for WS=50, and 0.10% for WS=150. We find that enlarging the WS reduces the DR further so that it is close to or equal to 0% for all network attacks. For example, attacks with fewer packets, such as the brute-force attack with 1,507 packets, can fully evade the IF NIDS when the window size is enlarged from three to 50.

For other attacks against the IF NIDS, e.g., the active wiretap attack, which has 500,000 network packets, the increase in window size has to be even larger. Here, increasing the window size to 150 results in almost perfect evasion, with a DR of 0.79%. Similar behavior is observed for the other NIDSs. For example, reshaping the SSDP flood attack (1,500,000 packets) with a window size of 150 against the AE NIDS results in a 100% success rate (0% DR). Overall, we find that all attacks are undetected when the window size is set at 150. Given the fact that an attacker usually does not know which NIDS is used by the target network, it thus makes sense to use an LSTM window size of 150.

Table 4: NIDS DR with increasing LSTM window size (in %).

NIDS/Attack Type	Before	After↓		
		WS=3	WS=50	WS=150
AE/Active Wiretap	98.03	0.00	0.79	0.00
AE/SSDP Flood	71.94	10.22	0.00	0.00
AE/Brute-Force	25.50	0.00	0.00	0.00
KitNET/Active Wiretap	90.02	0.00	0.00	0.20
KitNET/SSDP Flood	99.97	0.00	0.00	0.00
KitNET/Brute-Force	22.79	0.00	0.00	0.00
IF/Active Wiretap	69.73	7.29	7.09	0.73
IF/SSDP Flood	99.94	0.01	0.01	0.00
IF/Brute-Force	87.94	9.34	0.00	0.00
<b>Average</b>	<b>73.98</b>	<b>2.98</b>	<b>0.88</b>	<b>0.10</b>

**Impact Summary:** Having optimized TANTRA, all eight attacks remain undetected against all three NIDSs, decreasing the average DR to 0.01%. Compared to prior studies, this demonstrates an up to 28.18% improvement and a 12.89% average improvement (Table 3). Furthermore, the attacks functionality is preserved as only attack network traffic timestamps are reshaped.

### 6.3 End-to-End Evaluation

In the previous section, we examined the effectiveness of TANTRA when reshaping all of the attack network traffic and evaluated the attack against an NIDS; this enables us to perform a comparison to related work. One of TANTRA's key strengths, however, is its end-to-end capability, which allows the attacks to be done using an active target network connection, as would be the case for a real-world scenario. TANTRA faces a more difficult challenge in this setup, as it has to reshape attack network traffic step by step, according to the target network's delays and responses in real time.

To perform this evaluation, we (1) set up a target network with a KitNET NIDS, (2) train the timestamp generating LSTM on benign network traffic, (3) reshape the MITM malicious attack traffic using the trained LSTM, and (4) evaluate TANTRA against the NIDS. We then set up a proxy using the Scapy [3] library to send each packet with the specific timestamp generated by the LSTM model. In this way, network traffic is sent to the target network while adjusting the timestamps in real time.

Figure 6 shows KitNET's anomaly scores when processing TANTRA's reshaped MITM network traffic. The results indicate that the reshaping results in full evasion throughout the active connection.

One interesting observation is that after the first step, from the 348,000 packet onward, the anomaly score increases slightly. This could be due to additional network delay and

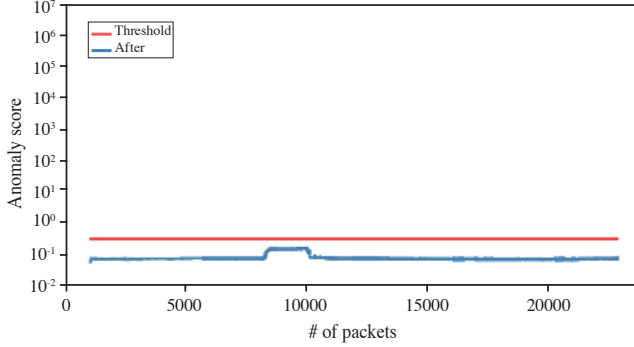


Figure 6: KitNET NIDS anomaly scores for MITM attack in end-to-end scenario.

response activity. This may illustrate the relevance of an independent end-to-end evaluation for validating the performance under real-world conditions. Despite the slight increases in the anomaly score, the attack’s scores remain below the NIDS detection threshold. Therefore, we consider TANTRA the first viable approach capable of overcoming the challenges of an end-to-end scenario, while remaining fully undetected and retaining the functionality of the attack network packet contents.

In future work, this scenario can be evaluated for other types of intrusion attacks. Furthermore, other target networks can be explored, enabling different types of proxy to be tested for intrusion. In addition, other types of NIDSs can be integrated. In addition, by optimizing the window size, TANTRA can be fine-tuned for specific attack, network, and NIDS characteristics, to obtain the greatest impact. Each of these settings can be evaluated using the general approach presented in this paper.

**Impact Summary:** While additional difficulty is imposed by target network delays and responses, by using TANTRA, attacks can evade detection throughout the entire intrusion process. The observed anomaly scores on KitNET show the impact of network responses, which resulted in a spike in the score. However, TANTRA’s ability to integrate such responses in the reshaping process resulted in a 0% DR for the evaluated attack.

## 6.4 Mitigating TANTRA

To identify leading directions for mitigating TANTRA, we propose a new way of training an NIDS. We use the methodology presented in Section 4.5, but instead of training on benign behavior for the detection of unknown anomalies, we train on both benign traffic and reshaped malicious traffic for the purpose of distinguishing between them. The suggested mitigation is evaluated on three supervised-learning ML-based NIDSs with the following settings:

- **Logistic Regression (LR).** We use an logistic regression NIDS following scikit-learn [27] default hyperparameter setting, where we do not execute normalization (`normalize=False`), and we do not force the coefficients to be positive (`positive=False`). Furthermore, maximum iterations of convergence are set to 100.
- **Gaussian Naive Bayes (NB).** The Gaussian Naive Bayes NIDS also follows the default settings by the scikit-learn, where the portion of the largest variance of all features that is added to variances for calculation stability is  $1e-9$  (`var_smoothing=1e-09`).
- **Random Forest (RF).** For the Random Forest NIDS, we use 100 estimators and the Gini criterion to measure the impurity of the estimator splits. We do not define a maximum depth and set the minimum sample splits to 2. Furthermore, bootstrap samples are used to build the trees.

The NIDSs are evaluated as described in Section 6.2. We use both small and large packet size attacks, namely active wiretap, MITM, fuzzing, and SSDP flood attacks. We evaluate the NIDSs using cross-validation. Benign traffic, along with three reshaped attacks, is used for training. Benign and reshaped malicious traffic is trained in equal proportion. The remaining fourth attack is used for testing.

The results presented in Table 5 show that the defense technique achieves a total average DR of 45.93% at the cost of a 20.97% false positive rate (FPR). Both, the NB NIDS and RF NIDS are able to detect almost all reshaped malicious traffic for the MITM and SSDP flood attacks. However, the NB NIDS is the only one that achieves such performance with a 0% FPR. In contrast, the RF NIDS has an FPR of 17.73%.

The least effective NIDS seems to be the LR NIDS whose average FPR is 6.53 percentage points higher than the average DR. A possible explanation for such behavior may be that the NB and RF NIDSs focus on dealing with continuous data, which may help find connections between different features. On the other hand, the LR NIDS aims to build a linear function to differentiate between classes, which may be more difficult in our case.

Interestingly, none of the NIDSs were able to detect the active wiretap or fuzzing attacks after TANTRA was used. One reason for that may be that both attacks are designed for execution on a router, while MITM and SSDP flood attacks are designed for execution inside a connection in the target network, as shown by Mirsky et al. [23]. Therefore, the attack design differentiates the four attacks, which may cause the NIDS training process to be less effective on two of them.

Overall, the defense technique seems to be highly effective for some attacks, while other attacks remain undetected. While not yet fully effective, the proposed mitigation provides a promising direction for addressing TANTRA.

**Challenges and Discussion.** To further mitigate TANTRA, general and versatile defense techniques must be identified. So far, changing the training paradigm has shown promising results for two out of the four evaluated attacks. One of the main challenges is to train NIDSs so that they are able to detect any kind of reshaped malicious traffic. Our proposed defense technique uses reshaped malicious traffic from TANTRA’s methodology to train the NIDS specifically for such traffic. However, as evasion attacks evolve, the trained NIDS may become outdated and unable to cope with the unfamiliar new reshaping methods. Hence, there is a need for a general defense technique that can detect reshaped malicious traffic, even when the method for reshaping is unknown to the defender.

The second challenge that mitigation faces has to do with the NIDS used to perform detection. In this mitigation evaluation, different NIDSs than those used when evaluating our attack. The reason is the autoencoder’s compatibility with non-binary decisions. Therefore, two different classes with four decision outcomes would require an individual AE, which then could be combined in an ensemble, confirming or disregarding one another; we propose pursuing this direction in future research.

Finally, the identified mitigation directions could be integrated into existing NIDSs which are trained solely on benign traffic. In such a case, the detection systems could work in a combined fashion in which the first one provides an anomaly score, and the second one is used to predict whether or not the input is a malicious or benign sample. However, the optimal scenario would require training the NIDSs, which are aware of reshaped malicious traffic, against all variants of evasion and maintaining such detection systems with the latest attacks, similar to traditional malware classifiers such as antivirus software.

## 7 Related Work

In this section, we present related work and compare the methods proposed in those studies to TANTRA, in terms of the preparation required before the attack is executed, the attack’s execution, the ability to remain undetected, the impact of the attack and other ways of mitigation.

The majority of related studies have focused on reshaping the feature space. Lin et al. [20] generated adversarial attacks using a unique type of generative adversarial network (GAN), namely IDSGAN. Yang et al. used internal information about the detection system to generate adversarial examples. Wang et al. [37] used a multi-layer perceptron (MLP), which was trained using the original training dataset, to generate adversarial examples using both the FGSM and JSMA adversarial attack methods. Warzyński and Kołaczek [38] attempted to compromise an NIDS based on a surrogate model. While Clement et al. [10] assumed that the adversary has direct knowledge of the target DL NIDS, allowing them to generate

Table 5: NIDS DR following novel defense technique.

NIDS	Attack Type	DR↓	FPR↓
Logistic Regression	Active Wiretap	0.00	50.70
	MITM	100.00	33.60
	Fuzzing	0.00	83.00
	SSDP Flood	54.60	13.40
	Average	38.65	45.18
Gaussian Naive Bayes	Active Wiretap	0.00	0.00
	MITM	100.00	0.00
	Fuzzing	0.00	0.00
	SSDP Flood	100.00	0.00
	Average	50.00	0.00
Random Forest	Active Wiretap	0.00	0.00
	MITM	100.00	39.90
	Fuzzing	0.00	0.00
	SSDP Flood	96.50	31.00
	Average	49.13	17.73
<b>Total Average</b>		<b>45.93</b>	<b>20.97</b>

input for the deep learning model directly. The work closest to ours is the study performed by Han et al. [16]. The authors did not use the feature space; instead, they proposed an evasion attack to bypass an NIDS, by combining the use of NIDS behavior analysis and dummy packets to reshape the traffic.

**Preparation.** Most related work assumed to have access to the targeted system’s feature extraction process [10, 20, 30, 37–39]. With regard to the preparation required for the attack, the approaches presented in other papers mainly relied on training a separate DL model that is used to perform an adversarial attack. In other cases, there was a need to obtain access to the targeted classifier’s input and evaluate its output [10, 30, 39]. While capable of bypassing an NIDS, the required knowledge makes these attacks impractical.

Gaining access to NIDS behavior relies on threat models which are both difficult to build and require information and knowledge that is difficult to obtain, which questions the attacks’ value when compared to the benefits of evading the NIDS. Hat et al. trained a GAN to produce reshaped traffic by integrating dummy packets into the attack network traffic. To do so, all of the traffic must be generated [16]. This approach neglects the issue of integrating target network delays or responses, which makes it inapplicable in an end-to-end scenario under the proposed setup. TANTRA’s preparation consists of training an LSTM on benign traffic obtained from the target network. In contrast to related studies, knowledge about the NIDS is not needed, meaning that TANTRA is a black-box evasion attack.

**Execution & Performance.** When comparing execution and performance, it is crucial to distinguish between white-box and black-box evasion attacks. White-box evasion attacks assume barely obtainable knowledge which may increase their performance. This performance, however, is rarely achievable in a real-world setting. Black-box evasion attacks are better suited for direct comparison to TANTRA.

In the case of white-box attacks, Clement et al. achieved 100% success for integrity attacks, while Lin et al. achieved similar results, as they decreased the DR to at most 1.56%. Rigaki et al. were able to decrease the targeted NIDS’ accuracy by 0.25%, on average, while Wang et al. were able to decrease the DR from an average of 56.128% to 20.0% for most combinations of adversarial and network attacks.

For black-box attacks, Han et al. used similar to our setup, the Mirai, fuzzing, SSDP flood, and brute-force attacks and evaluated the attacks after reshaping on a KitNET NIDS. Thereby the authors’ approach decreased the average detection rate to 12.89%. In contrast, the use of TANTRA decreased the DR to 0%. This makes it possible to fully evade those attacks against a KitNET NIDS; we also validate our method on two other state-of-the-art NIDSs (AE and IF NIDSs). Furthermore, TANTRA is able to obtain similar performance for the other attacks evaluated in this work.

This comparison shows that TANTRA outperforms the related black-box evasion attacks and even outperforms white-box evasion attacks which have far greater knowledge on the NIDS.

Execution-wise, in comparison to related work, TANTRA simply reshapes attack network traffic using an LSTM and the timestamp attribute. This may limit the flexibility in reshaping the traffic, however it guarantees functional attack network traffic. Moreover, TANTRA is able to interact with the target network in each step of the attack in real-time, allowing it to be effective in more difficult settings.

**Impact.** The impact of bypassing an NIDS can be great. For example, it can allow an attacker to take control of a targeted network, listen to network traffic, or deny service. The attacks presented in other studies that use the feature space rarely have an impact on the target network. First, such attacks require knowledge on the NIDS that is very difficult to obtain. Second, even if the knowledge is obtained, the attack is likely to lose its functionality when the reshaped feature space is transformed back to the problem space (Section 4.1). These studies may serve as the basis for the development of more impactful attacks if a means for such reconstruction is developed in the future.

On the other hand, the black-box evasion attacks presented in Han et al. [16] may have impact. Regarding the attack presented in that paper, there is a risk that the attack will be inexecutable given the use of dummy packages or be detected because of the inability to integrate network responses or delays. But there is still a chance that the approach will be

successful for some other cases in a real-world setting. In contrast, TANTRA’s effectiveness was demonstrated in a full end-to-end scenario with an active network connection, where it was 100% successful in evading detection by the KitNET NIDS for various attacks.

**Mitigation.** Some mitigation has been proposed by prior research. One of the proposed defense technique is adversarial feature reduction [16]. Here, dimensions of features are reduced to those which provide the NIDS the greatest insight into distinguishing between malicious and benign packets. While this may be effective against the attack presented in their paper, it is not the case for TANTRA. Since TANTRA’s approach focuses on emulating benign behavior, adversarial feature reduction will have difficulties in defending against and difficulties to detect the reshaped malicious traffic of TANTRA.

Other studies trained networks on adversarial robustness, exposing the NIDS to adversarial attacks during training, in order to guide anomaly detection [2]. Similar behavior is observed in the emerging area of out-of-distribution detection. In this case, the system is exposed to outliers, called outlier exposure [17], which is then integrated into the training procedure using a special loss function. As a result, the system behaves differently to the outliers which are detected by a drop in prediction accuracy or overall prediction entropy.

Overall, with the help of the defense technique presented in this paper, we identified that exposing the system to reshaped malicious traffic helps improve detection for some attacks. However, the proposed defense technique still requires optimization for being considered a generic approach to evading evasion attacks.

## 8 Conclusion

In this work, we presented TANTRA, the first end-to-end evasion attack that can evade all state-of-the-art NIDSs with a 99.99% success rate. We showed that by simply reshaping attack network traffic using the timestamp attribute, TANTRA outperforms related work’s evasion attacks by up to 28.18%. Furthermore, the attack was capable of successfully evading the NIDS when executed in a black-box setting. When evaluated under an active network connection, TANTRA bypassed KitNET, one of the most advanced NIDSs, and is the first evasion attack that does not affect the attack packets’ content. Therefore, the attack’s functionality is retained, allowing the attack to impact the target network.

We also presented a promising mitigation direction that involves training a NIDS with both benign and reshaped traffic. While not yet fully effective against all types of attacks, this defense technique shows to be a promising direction and is worth pursuing in future research, as a means of defending against the TANTRA threat presented in this paper.

## References

- [1] Intriguing properties of adversarial ml attacks in the problem space. *arXiv preprint arXiv:1911.02142* (2019).
- [2] AL-DUJAILI, A., HUANG, A., HEMBERG, E., AND O'REILLY, U.-M. Adversarial deep learning for robust detection of binary encoded malware, 2018.
- [3] BIONDI, P. Scapy documentation (!), 2010.
- [4] BUCZAK, A. L., AND GUVEN, E. A survey of data mining and machine learning methods for cyber security intrusion detection. *IEEE Communications surveys & tutorials* 18, 2 (2015), 1153–1176.
- [5] CAO, J., LI, Z., AND LI, J. Financial time series forecasting model based on ceemdan and lstm. *Physica A: Statistical Mechanics and its Applications* 519 (2019), 127–139.
- [6] CARLINI, N., AND WAGNER, D. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)* (2017), IEEE, pp. 39–57.
- [7] CHALAPATHY, R., AND CHAWLA, S. Deep learning for anomaly detection: A survey. *arXiv preprint arXiv:1901.03407* (2019).
- [8] CHEN, P.-Y., SHARMA, Y., ZHANG, H., YI, J., AND HSIEH, C.-J. Ead: elastic-net attacks to deep neural networks via adversarial examples. *arXiv preprint arXiv:1709.04114* (2017).
- [9] CHOLLET, F. ET AL. Keras, 2015.
- [10] CLEMENTS, J., YANG, Y., SHARMA, A., HU, H., AND LAO, Y. Rallying adversarial techniques against deep learning for network security. *arXiv preprint arXiv:1903.11688* (2019).
- [11] DING, Z., AND FEI, M. An anomaly detection approach based on isolation forest algorithm for streaming data using sliding window. *IFAC Proceedings Volumes* 46, 20 (2013), 12–17.
- [12] FOLORUNSO, O., AYO, F. E., AND BABALOLA, Y. Ca-nids: A network intrusion detection system using combinatorial algorithm approach. *Journal of Information Privacy and Security* 12, 4 (2016), 181–196.
- [13] GANDHI, A., AND JAIN, S. Adversarial perturbations fool deepfake detectors. *arXiv preprint arXiv:2003.10596* (2020).
- [14] GOODFELLOW, I. J., SHLENS, J., AND SZEGEDY, C. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572* (2014).
- [15] GÜERA, D., AND DELP, E. J. Deepfake video detection using recurrent neural networks. In *2018 15th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)* (2018), IEEE, pp. 1–6.
- [16] HAN, D., WANG, Z., ZHONG, Y., CHEN, W., YANG, J., LU, S., SHI, X., AND YIN, X. Practical traffic-space adversarial attacks on learning-based nids. *arXiv preprint arXiv:2005.07519* (2020).
- [17] HENDRYCKS, D., MAZEIKA, M., AND DIETTERICH, T. Deep anomaly detection with outlier exposure. In *International Conference on Learning Representations* (2019).
- [18] KWON, D., KIM, H., KIM, J., SUH, S. C., KIM, I., AND KIM, K. J. A survey of deep learning-based network anomaly detection. *Cluster Computing* (2019), 1–13.
- [19] LI, D., KOTANI, D., AND OKABE, Y. Improving attack detection performance in nids using gan. In *2020 IEEE 44th Annual Computers, Software, and Applications Conference (COMPSAC)* (2020), IEEE, pp. 817–825.
- [20] LIN, Z., SHI, Y., AND XUE, Z. Idsgan: Generative adversarial networks for attack generation against intrusion detection. *arXiv preprint arXiv:1809.02077* (2018).
- [21] LIU, F. T., TING, K. M., AND ZHOU, Z.-H. Isolation forest. In *2008 Eighth IEEE International Conference on Data Mining* (2008), IEEE, pp. 413–422.
- [22] LIU, F. T., TING, K. M., AND ZHOU, Z.-H. Isolation-based anomaly detection. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 6, 1 (2012), 1–39.
- [23] MIRSKY, Y., DOITSHMAN, T., ELOVICI, Y., AND SHABTAI, A. Kitsune: an ensemble of autoencoders for online network intrusion detection. *arXiv preprint arXiv:1802.09089* (2018).
- [24] MIRSKY, Y., AND LEE, W. The creation and detection of deepfakes: A survey. *arXiv preprint arXiv:2004.11138* (2020).
- [25] OLAH, C. Understanding lstm networks.
- [26] PAPERNOT, N., MCDANIEL, P., JHA, S., FREDRIKSON, M., CELIK, Z. B., AND SWAMI, A. The limitations of deep learning in adversarial settings. In *2016 IEEE European symposium on security and privacy (EuroS&P)* (2016), IEEE, pp. 372–387.
- [27] PEDREGOSA, F., VAROQUAUX, G., GRAMFORT, A., MICHEL, V., THIRION, B., GRISEL, O., BLONDEL, M., PRETTENHOFER, P., WEISS, R., DUBOURG, V., VANDERPLAS, J., PASSOS, A., COURNAPEAU, D., BRUCHER, M., PERROT, M., AND DUCHESNAY, E. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* 12 (2011), 2825–2830.
- [28] PUGGINI, L., AND MCLOONE, S. An enhanced variable selection and isolation forest based methodology for anomaly detection with oes data. *Engineering Applications of Artificial Intelligence* 67 (2018), 126–135.
- [29] RAGHUNATH, B. R., AND MAHADEO, S. N. Network intrusion detection system (nids). In *2008 First International Conference on Emerging Trends in Engineering and Technology* (2008), pp. 1272–1277.
- [30] RIGAKI, M. Adversarial deep learning against intrusion detection classifiers, 2017.
- [31] ROONDIWALA, M., PATEL, H., AND VARMA, S. Predicting stock prices using lstm. *International Journal of Science and Research (IJSR)* 6, 4 (2017), 1754–1756.
- [32] SCHNEIDER, P., AND BÖTTINGER, K. High-performance unsupervised anomaly detection for cyber-physical system networks. In *Proceedings of the 2018 Workshop on Cyber-Physical Systems Security and Privacy* (2018), pp. 1–12.
- [33] SHARAFALDIN, I., LASHKARI, A. H., AND GHORBANI, A. A. Toward generating a new intrusion detection dataset and intrusion traffic characterization. In *ICISSP* (2018), pp. 108–116.
- [34] SYMANTEC, Feb 2019.
- [35] SZEGEDY, C., ZAREMBA, W., SUTSKEVER, I., BRUNA, J., ERHAN, D., GOODFELLOW, I., AND FERGUS, R. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199* (2013).
- [36] WANG, P., CHEN, X., YE, F., AND SUN, Z. A survey of techniques for mobile service encrypted traffic classification using deep learning. *IEEE Access* 7 (2019), 54024–54033.
- [37] WANG, Z. Deep learning-based intrusion detection with adversaries. *IEEE Access* 6 (2018), 38367–38384.
- [38] WARZYŃSKI, A., AND KOLACZEK, G. Intrusion detection systems vulnerability on adversarial examples. In *2018 Innovations in Intelligent Systems and Applications (INISTA)* (2018), IEEE, pp. 1–4.
- [39] YANG, K., LIU, J., ZHANG, C., AND FANG, Y. Adversarial examples against the deep learning based network intrusion detection systems. In *MILCOM 2018-2018 IEEE Military Communications Conference (MILCOM)* (2018), IEEE, pp. 559–564.