

Apollon: A Robust Defence System against Adversarial Machine Learning Attacks in Intrusion Detection Systems

Antonio Paya^{a,*}, Sergio Arroni^a, Vicente García-Díaz^a and Alberto Gómez^b

^aDepartment of Computer Science, University of Oviedo, Science Faculty, Oviedo, Spain

^bDepartment of Business Administration, University of Oviedo, Gijón, Spain

ARTICLE INFO

Keywords:

Adversarial Machine Learning
Intrusion Detection Systems
Artificial Intelligence
Cybersecurity
Multi-Armed Bandits

ABSTRACT

The rise of Adversarial Machine Learning (AML) attacks is presenting a significant challenge to Intrusion Detection Systems (IDS) and their ability to detect threats. To address this issue, we introduce *Apollon*, a novel defence system that can protect IDS against AML attacks. *Apollon* utilizes a diverse set of classifiers to identify intrusions and employs *Multi-Armed Bandits (MAB)* with *Thompson sampling* to dynamically select the optimal classifier or ensemble of classifiers for each input. This approach enables *Apollon* to prevent attackers from learning the IDS behaviour and generating adversarial examples that can evade the IDS detection. We evaluate *Apollon* on several of the most popular and recent datasets, and show that it can successfully detect attacks without compromising its performance on traditional network traffic. Our results suggest that *Apollon* is a robust defence system against AML attacks in IDS.

1. Introduction

As computer networks continue to grow in size and complexity, the need for effective security measures becomes increasingly critical. Intrusion Detection Systems (IDS) have been developed to address this challenge, providing a means of monitoring network traffic and identifying potential security threats. These systems can analyze network traffic and identify potential security threats such as malware, network intrusions, and denial of service attacks. However, the increasing complexity and diversity of network traffic have made it difficult to accurately classify network traffic using traditional rule-based IDS systems [1].

To overcome these limitations, Machine Learning (ML) techniques have been widely adopted in IDS for network traffic classification. These techniques leverage the power of statistical models and algorithms to automatically learn and detect anomalous network traffic patterns, which are indicative of security threats.

ML-based IDS systems offer several advantages over traditional rule-based systems, including higher accuracy, better scalability, and more robustness to evolving network threats [2, 3]. However, they also pose new challenges, particularly in terms of security. One of the main challenges is the susceptibility of ML models to Adversarial Machine Learning (AML) attacks [4].

AML attacks are a type of cyber-attack that aims to manipulate ML models by feeding them carefully crafted input data, called adversarial examples. These examples are designed to cause misclassification or incorrect predictions,

which can be exploited by attackers to bypass the security measures of IDS systems [5, 6, 7].

Attackers create adversarial examples by utilizing information obtained from the targeted IDS, including its responses to specific inputs. This information is used to train a model capable of generating adversarial traffic that remains undetectable by the IDS classifier.

As ML-based IDS systems become more prevalent, the threat of Adversarial Machine Learning (AML) attacks becomes more significant [8]. Therefore, it is crucial to develop effective defense mechanisms to mitigate the impact of these attacks and ensure the reliability and robustness of IDS systems.

In this paper, we propose a new robust defense system against Adversarial Machine Learning attacks on Intrusion Detection Systems called *Apollon*. *Apollon* serves to safeguard IDS from attackers by obstructing their ability to generate adversarial traffic through learning from the behavior of the IDS.

Apollon utilizes a diverse range of classifiers to detect intrusions and employs a *Multi-Armed Bandits (MAB)* with *Thompson sampling* to select the optimal classifier or a combination of classifiers in real-time for each input, enabling it to achieve this objective without compromising its performance on traditional network traffic.

In this way, *Apollon* can prevent attackers from learning the behavior of the IDS in realistic training times, adding a layer of uncertainty to the IDS behavior that makes it more difficult for attackers to detect the IDS behavior and generate adversarial traffic.

The structure of the paper is as follows. Section 2 provides an overview of the main concepts and techniques used in this paper. Section 3 discusses the related work in the field of AML attacks and IDS classifiers. Section 4 presents the proposed defense system, *Apollon*. Section 5 presents the experimental evaluation of *Apollon*. Finally, Section 6 concludes the paper and discusses future work.

*Corresponding author

antonio paya@outlook.com (A. Paya); sergioarroni@outlook.com (S. Arroni); garciavicente@uniovi.es (V. García-Díaz); albertogomez@uniovi.es (A. Gómez)

ORCID(S): 0000-0003-3733-3805 (A. Paya); 0000-0002-4907-8576 (S. Arroni); 0000-0003-2037-8548 (V. García-Díaz); 0000-0003-3570-4043 (A. Gómez)

2. Background

In this section, we will provide an overview of the background knowledge required to understand the rest of this work. This includes an introduction to the concepts of Adversarial Machine Learning and adversarial examples, as well as an introduction to the concepts of Intrusion Detection Systems and the challenges they face.

2.1. Intrusion Detection Systems

Intrusion Detection Systems (IDS) are security tools designed to identify and prevent unauthorized access, misuse, and malicious activities in computer networks [9]. IDS play a critical role in protecting networks from various types of cyber threats, including viruses, malware, and intrusions. IDS operate by monitoring network traffic and analyzing it for suspicious behavior or patterns. There are two main types of IDS: Network-based Intrusion Detection Systems (NIDS) and Host-based Intrusion Detection Systems (HIDS) [10].

NIDS monitors network traffic and analyzes packets to identify potential security threats. NIDS can be deployed at different points within the network, such as at the perimeter, within the LAN, or at critical points within the network. NIDS can detect a wide range of attacks, including port scans, denial-of-service attacks, and data exfiltration.

HIDS, on the other hand, monitors the activity on individual hosts, such as servers or workstations. HIDS can detect attacks that may not be visible to NIDS, such as attacks that occur within encrypted traffic or attacks that originate from within the network.

From this point on, and to facilitate the understanding of the document, we will refer to Network Intrusion Detection Systems as Intrusion Detection Systems.

Intrusion Detection Systems are an essential component of a comprehensive network security strategy. They provide an additional layer of protection beyond firewalls, antivirus software, and other security tools. By detecting and alerting administrators to potential security threats, IDS can help organizations respond quickly and effectively to cyber attacks.

Machine Learning (ML) has emerged as a powerful technique for improving the accuracy and effectiveness of Intrusion Detection Systems [2, 3, 1]. ML algorithms can be used to analyze large volumes of network data and identify patterns that may be indicative of security threats. ML-based IDS can learn from past network activity to identify and flag potential security threats in real-time, even when the attacks are novel or previously unseen. ML-based IDS can also adapt and improve over time as they learn from new data and feedback from security analysts.

2.2. Adversarial Machine Learning

Adversarial Machine Learning (AML) attacks refer to a set of techniques used to undermine the accuracy, integrity, or security of machine learning (ML) models [4]. AML attacks can be launched by malicious actors with different objectives, such as stealing sensitive information, manipulating decision-making processes, or compromising the confidentiality and privacy of ML systems.

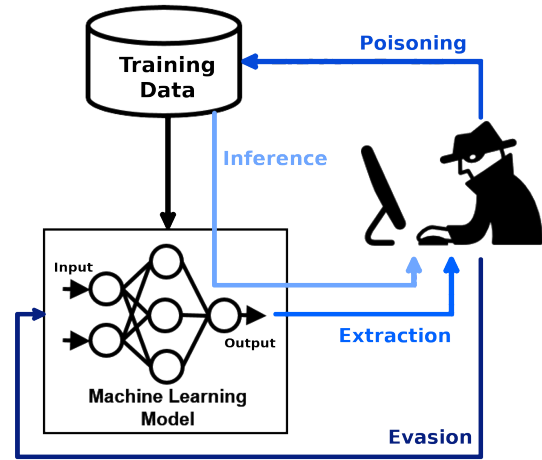


Figure 1: Adversarial Machine Learning (AML) attacks by ART [11]

AML attacks can be launched against a wide range of ML models, including deep neural networks, support vector machines, decision trees, and others. The success of an AML attack depends on various factors, such as the type and quality of the target ML model, the sophistication of the attack technique, and the attacker's level of knowledge and resources. According to the taxonomy of the attack, they can be classified into evasion attacks, poisoning attacks, extraction attacks and inference attacks[12].

2.2.1. Evasion Attacks

Evasion attacks in AML refer to a type of attack where the attacker manipulates the input data in a way that the ML model will misclassify it, without changing the underlying characteristics of the data. Evasion attacks are typically launched against classification models, such as those used for image recognition or spam detection, and they can be crafted using various techniques, including gradient-based methods, evolutionary algorithms, or grey/black-box attacks. The goal of an evasion attack is to create an adversarial example, i.e., a modified version of the original input data that is similar to the original but is misclassified by the ML model. Evasion attacks pose a significant threat to the security and robustness of ML systems, especially in domains such as malware detection, Intrusion Detection Systems and fraud detection, where accurate classification is critical.

2.2.2. Poisoning Attacks

Data poisoning attacks in AML involve manipulating the training data of an ML model to introduce biases or to cause it to learn incorrect patterns. Poisoning attacks can be launched at different stages of the ML pipeline, including data collection, preprocessing, and training. The goal of a poisoning attack is to compromise the integrity and accuracy of the ML model by introducing malicious data into the training dataset, which can cause the model to learn incorrect patterns and make incorrect predictions. Poisoning attacks

can be launched in a variety of ways, such as by injecting adversarial examples into the training data, manipulating the distribution of the training data, or introducing outliers into the dataset.

2.2.3. Extraction Attacks

Model extraction attacks in AML refer to a type of attack where the attacker aims to extract the details of an ML model without direct access to it. This is achieved by leveraging the output of the target ML model to infer the underlying structure, architecture, or parameters of the model. Model extraction attacks can be launched through different channels, such as querying the model with carefully crafted inputs or by observing its behavior in response to various inputs. The goal of a model extraction attack is to steal the target model's intellectual property or use it for malicious purposes such as deploying counterfeit models, stealing sensitive data, or reverse engineering proprietary algorithms. Model extraction attacks can be particularly effective against black-box models where the attacker does not have access to the model's internal structure or parameters.

2.2.4. Inference Attacks

AML inference attacks involve an attacker attempting to glean confidential information about the input data utilized by the ML model by scrutinizing the model's output. Inference attacks can be launched against a wide range of ML models, including deep neural networks, decision trees, and support vector machines, among others. The goal of an inference attack is to obtain access to private or confidential information about the input data, such as personal characteristics, financial transactions, or medical records, without having direct access to the data itself. Inference attacks can be launched through different channels, such as analyzing the output distribution of the model, measuring its response time to different inputs, or exploiting the model's decision boundaries.

2.3. Multi-Armed Bandits (MAB)

The *Multi-Armed Bandits (MAB)* is a classic problem in probability theory and Machine Learning, where an agent has to allocate a limited set of resources among competing choices that have uncertain rewards [13]. The agent faces a trade-off between exploiting the choices that have the highest expected rewards based on the current information, and exploring new choices that may yield higher rewards in the future.

The *MAB* problem has many practical applications in various domains, such as clinical trials, adaptive routing, financial portfolio design, and online advertising. Several algorithms have been proposed to solve the *MAB* problem, such as optimistic initialization [14], upper confidence bound (UCB) [15], and Thompson sampling [16]. These algorithms differ in how they balance exploration and exploitation, and how they estimate the expected rewards of each choice.

Thompson sampling is a Bayesian approach that maintains a probability distribution over the unknown reward

distributions of each choice, and chooses actions based on sampling from these distributions. Specifically, at each timestep, Thompson sampling samples a reward from each distribution, chooses the action associated with the highest sampled reward, and updates its beliefs about the reward distributions based on the observed reward. This approach has been shown to be effective in many applications, and has a strong theoretical justification in terms of minimizing regret.

Thompson sampling has gained popularity in recent years due to its ability to balance exploration and exploitation in a principled way [17]. By sampling from the probability distributions over the reward distributions, Thompson sampling encourages exploration of all choices while still favoring choices with higher expected rewards. Additionally, the Bayesian framework allows for the incorporation of prior knowledge about the reward distributions, which can be especially useful in scenarios with limited data.

The *MAB* problem is intricately connected to the realm of Reinforcement Learning (RL). In RL, an intelligent agent endeavors to acquire knowledge and develop a strategy, known as a policy, that maximizes its total rewards throughout its interaction with an environment. Over the past few years, RL has exhibited remarkable success in diverse domains. Notably, it has found significant utility in the domain of demand forecasting [18, 19].

In the context of our work, we use the *MAB* algorithm to select the best (IDS) classifier for each network traffic request. This is similar to how *MAB* is used in demand forecasting to select the most optimal forecasting model or determine the best hyperparameters for a given forecasting model. By using *MAB*, we can balance the trade-off between exploiting the classifiers that have the highest expected accuracy based on the current information, and exploring new classifiers that may yield higher accuracy in the future.

3. Related Work

This section will present a summary of the distinct datasets for IDS, along with their corresponding classifiers and performance metrics. Our proposed approach will make use of these classifiers. Furthermore, we will explore the typical types of AML attacks used in this domain.

3.1. Intrusion Detection Systems datasets

IDS datasets play a crucial role in assessing and gauging the performance of IDSs. These datasets contain labeled instances of regular and anomalous network traffic that are used to train and assess the precision and efficiency of IDSs. A range of datasets with diverse strengths and features are accessible. This section will examine some of the most prevalent datasets, highlighting their essential qualities and applications.

3.1.1. Discarded datasets

This project did not make use of several other IDS datasets, including Darpa 1998/99 [20], KDD 99 [21], and NSL-KDD [21]. These datasets are no longer commonly

employed for evaluating and benchmarking IDS due to their outdated nature. Created in the late 1990s and early 2000s, they do not accurately represent the current landscape of network threats and behaviors. KDD 99 dataset, in particular, has been criticized for its high false positive rate and lack of realism, thereby limiting its usefulness in assessing the performance of modern IDS [22, 23].

3.1.2. CIC-IDS-2017

A highly utilized IDS dataset in contemporary literature is the CIC-IDS2017 [24], which was developed by the Canadian Institute for Cybersecurity (CIC) in a simulated enterprise network environment, gathering network traffic data for five consecutive days. This dataset emulates the actions of 25 users and comprises nearly 80 significant attributes [25]. Notably, it has an 83% to 17% benign to malicious instance ratio, representing a significant portion of the dataset. The CIC-IDS2017 is considered an accurate depiction of normal traffic distribution in a network and can be utilized individually or combined with other datasets [26].

3.1.3. CSE-CIC-IDS-2018

The CSE-CIC-IDS2018 dataset was developed using AWS resources in a simulated enterprise network environment in 2018 [27]. It consists of data on seven distinct attack categories and comprises nearly 79 important features. With over 450 devices, including servers, computers, and other tools, this dataset is notably large and realistic [28]. It is akin to the CIC-IDS2017 dataset, analyzing bidirectional flow packet data, but with more significant features and greater comprehensiveness. Hence, it is widely used in the literature for assessing and benchmarking IDSs [28].

3.1.4. CIC-DDoS-2019

To address the lack of representation of all DDoS (Distributed Denial of Service) attack subtypes in existing datasets, the CIC-DDoS-2019 dataset was created [29]. Although the dataset includes simulated network traffic, it strives to present realistic benign data. It features 13 types of DDoS attacks and over 80 significant features. However, it is severely imbalanced, with 50,006,249 DDoS attack records and just 56,863 benign traffic records, making it challenging to train a model on both data types [25]. As a result, experts suggest using this dataset in conjunction with other datasets [26], such as CIC-IDS-2017 or CSE-CIC-IDS-2018, to train a more robust model.

3.2. Intrusion Detection Systems ML-classifiers

ML-classifiers have emerged as a promising alternative to traditional IDS for detecting network attacks. This is due to the limitations of traditional IDS in dealing with the complex and dynamic nature of cyber-attacks. In the current digital era, the number and sophistication of malware threats are constantly growing, posing a serious challenge to network security. Therefore, it is essential to have reliable and effective IDS systems in place to protect network systems from potential damage.

Several studies have been conducted to evaluate the performance of various ML-classifiers in detecting network attacks. These studies use datasets such as CIC-IDS-2017, CSE-CIC-IDS-2018, and CIC-DDoS-2019.

Table 1 provides a summary of the most commonly used ML-classifiers and their scores, including accuracy, F1 Score, and AUC, for each of the aforementioned datasets. The ML-classifiers used in these studies are *Logistic Regression (LR)* [36], *Fuzziness based Neural Networks (FNN)* [37], *Random Forests (RF)* [38], *Decision Trees (DT)* [39], *Robust transformer based Intrusion Detection System (RTIDS)* [33], and *Support Vector Machines (SVM)* [40].

These classifiers were selected due to their wide adoption and proven effectiveness in diverse Machine Learning applications. By examining the results presented in Table 1, researchers can gain a comprehensive understanding of the performance of these classifiers on the specific datasets, enabling informed decisions when choosing the most appropriate algorithm for their requirements. The detailed explanations of how each of the results were obtained for each classifier can be found in the papers listed in the *References* column. These papers provide comprehensive insights into the methodologies employed and offer further analysis of the performance of each ML-classifier on the respective datasets.

Accuracy, F1 Score, and AUC are three common metrics used to evaluate the performance of machine learning models. accuracy measures the proportion of correct predictions made by a model out of the total number of predictions and is defined as:

$$accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

where TP is the number of true positives, TN is the number of true negatives, FP is the number of false positives, and FN is the number of false negatives.

F1 Score is a weighted average of precision and the detection rate (DR) (also known as recall or sensitivity), where precision measures the proportion of true positives out of all predicted positives, and the detection rate measures the proportion of true positives out of all actual positives. The F1 Score is defined as:

$$F1Score = 2 \cdot \frac{precision \cdot DR}{precision + DR}$$

where precision is defined as:

$$precision = \frac{TP}{TP + FP}$$

and detection rate is defined as:

$$DR = \frac{TP}{TP + FN}$$

	CIC-IDS-2017			CSE-CIC-IDS-2018			CIC-DDoS-2019			
Classifiers	Accuracy	F1 Score	AUC	Accuracy	F1 Score	AUC	Accuracy	F1 Score	AUC	References
LR	92.96	90.87	91.50	87.96	88.99	81.54	91.72	87.27	90.23	[30, 31]
FNN	99.61	99.57	99.83	93.00	92.00	100.00	95.55	95.50	95.63	[32, 30, 33]
RF	99.79	99.78	99.98	92.00	94.00	100.00	99.86	99.78	99.82	[28, 32, 34, 3, 35, 30]
DT	99.62	99.57	99.56	88.00	91.00	100.00	99.87	99.78	99.80	[28, 30, 32, 3]
RTIDS	99.35	99.17	98.83	-	-	-	98.58	98.48	98.66	[33]
SVM	96.97	96.99	98.98	61.00	66.00	100.00	94.02	94.98	94.24	[28, 32, 3, 35, 33]

Table 1

Performance of the IDSs classifiers on the selected datasets.

The Area Under the ROC Curve (AUC) is a widely used performance metric in ML and binary classification tasks. It quantifies the discriminative power of a classification model by measuring the probability that the model will rank a randomly chosen positive instance higher than a randomly chosen negative instance. The ROC curve plots the true positive rate (DR) against the false positive rate (1-specificity) for various classification thresholds, where specificity is defined as:

$$specificity = \frac{TN}{TN + FP}$$

The AUC represents the integral of this curve and ranges from 0 to 1, where a value of 1 indicates a perfect classifier and a value of 0.5 suggests a random or ineffective classifier. Higher AUC values indicate better model performance in distinguishing between positive and negative instances. The AUC is a popular evaluation metric as it is robust to class imbalance and provides a concise summary of the model's overall performance.

Among the classifiers in Table 1, *Random Forest* [41] and *Decision Trees* [42] are found to be some of the most effective classifiers for detecting network attacks. *Random Forest* has gained popularity due to its ability to handle large datasets and its robust performance even when the data contains noise or missing values. *Decision Trees* are also preferred because of their simplicity and interpretability. They enable clear visualization of the decision-making process, making them useful for understanding the factors that contribute to the classification results.

Overall, these classifiers have demonstrated strong performance in the field of IDS and are frequently used by researchers and practitioners. Their effectiveness in detecting intrusions and classifying network traffic makes them valuable tools for maintaining the security and integrity of computer networks.

3.3. Adversarial Machine Learning attacks

ML-based IDSs can learn from data and adapt to new situations, unlike traditional systems that rely on predefined rules. However, ML-based IDSs also face new challenges from attackers who use Artificial Intelligence (AI) to craft sophisticated attacks that can fool or compromise ML models. One such threat comes from the use of Artificial Intelligence (AI) in the form of Adversarial Machine Learning (AML), where attackers use sophisticated techniques to manipulate or subvert ML models. These attacks are attractive to cyber attackers since they can be challenging to detect and prevent. Furthermore, as AI techniques gain popularity in cybersecurity, attackers are incentivized to develop more sophisticated adversarial attacks to evade detection.

Adversarial Machine Learning attacks can be classified as white-box attacks and grey/black-box attacks, depending on the level of knowledge the attacker possesses about the target model.

3.3.1. White-box attacks

White-box attacks are the most powerful kind of AML attacks because they let the attacker know everything about the IDS classifier and the training data. With this knowledge, the attacker can create very specific and complex attacks that can evade the system's defenses and achieve their goals. However, white-box attacks are also the most unrealistic kind of attack because they need the attacker to have a lot of knowledge about the system and its weaknesses, which is often impossible. In most situations, the attacker will only know some or little about the system, making white-box attacks impossible. Therefore, white-box attacks are very uncommon in reality and are usually only done in very focused and well-planned operations.

White-box attacks commonly used on IDS include the *Fast Gradient Sign Method (FGSM)* [43], *Deep-Fool* [44], *Carlini & Wagner attack (C&W)* [45], *Jacobian based Saliency Map Attack (JSMA)* [46], *Basic Iterative Method (BIM)* [47], and *Projected Gradient Descent (PGD)* [48].

3.3.2. Grey/Black-box attacks

In contrast to white-box attacks, grey/black-box attacks are more realistic because they do not require the attacker to have knowledge about the target model. However, these attacks are often less efficient compared to white-box attacks because the attacker's limited knowledge about the system restricts their ability to create highly targeted and sophisticated attacks. Instead, they have to depend on more general techniques that may be less effective in bypassing the system's defenses.

Generative Adversarial Networks (GANs) [49] are frequently used in black-box and grey-box attacks to produce adversarial examples. GANs are a type of Machine Learning model consisting of two neural networks: a generator network and a discriminator network. The generator network is trained to produce synthetic data that resembles the real data, while the discriminator network is trained to differentiate between real and synthetic data.

In the context of black-box and grey-box attacks, the generator network can generate adversarial examples that are specifically designed to evade the target system's defenses. The attacker may have access to the system's model scores or just a binary output indicating whether the input was accepted or rejected. This information can be utilized to guide the training of the generator network, enhancing its ability to produce effective adversarial examples.

Recent grey/black-box attacks on IDS include *attackGAN* [5], *DIGFuPAS* [8], *IDSGAN* [6], *VulnerGAN* [7], *ZOO attack* [50], *Boundary attack* [51] and the *HotSkipJump attack (HSJA)* [52]. *ZOO* is a score-based attack that estimates gradients to create adversarial traffic in grey/black-box settings. *Boundary attack* and *HSJA* are decision-based attacks that only use binary feedback to craft adversarial inputs. The *IDSGAN*, *attackGAN*, and *DIGFuPAS* are grey/blackbox attacks that employ *Wasserstein-GAN* to generate adversarial traffic. *Wasserstein-GAN (W-GAN)* [53] is a GAN variant that trains the generator network with a different objective function called the Wasserstein distance. The Wasserstein distance measures the distance between two probability distributions and has properties like smoothness and continuity. Some recent WGANs use the Gradient Penalty to enhance the training convergence.

3.4. Adversarial Machine Learning defenses

Due to the increasing number of AML attacks, researchers have developed several defense mechanisms to mitigate the impact of these attacks. In the IDS domain, these defenses can be classified into three categories [54]: *preprocessing defenses*, *adversarial training defenses* and *adversarial detection defenses*.

3.4.1. Preprocessing defenses

In order to mitigate the impact of adversarial perturbations, researchers have devised carefully planned preprocessing techniques. These preprocessing methods aim to reduce the vulnerability of machine learning models to adversarial attacks and enhance their robustness applying carefully planned transformations to the input data before it

is fed into the model. These transformations are designed to reduce the vulnerability of the model to adversarial perturbations.

One example of a preprocessing defense technique is *Stochastic Transformation-based Defenses* [55]. This technique involves applying random transformations to the input data, such as rotations, translations, and scaling, before feeding it into the model. By introducing randomness into the input data, the model becomes less susceptible to adversarial perturbations that are designed to exploit specific features of the input.

Another example of a preprocessing defense technique is *Gradient Masking* [56]. This technique involves modifying the gradients of the model during training to make it more difficult for an attacker to compute the gradients needed to generate adversarial examples. This is achieved by adding noise to the gradients or by clipping them to a certain range.

Nevertheless, sophisticated adversarial attacks have proven the inadequacy of these defense mechanisms. The primary flaws in these strategies lie in their ability to merely "confound" adversaries, rather than completely eliminating the presence of adversarial examples [57].

3.4.2. Adversarial training defenses

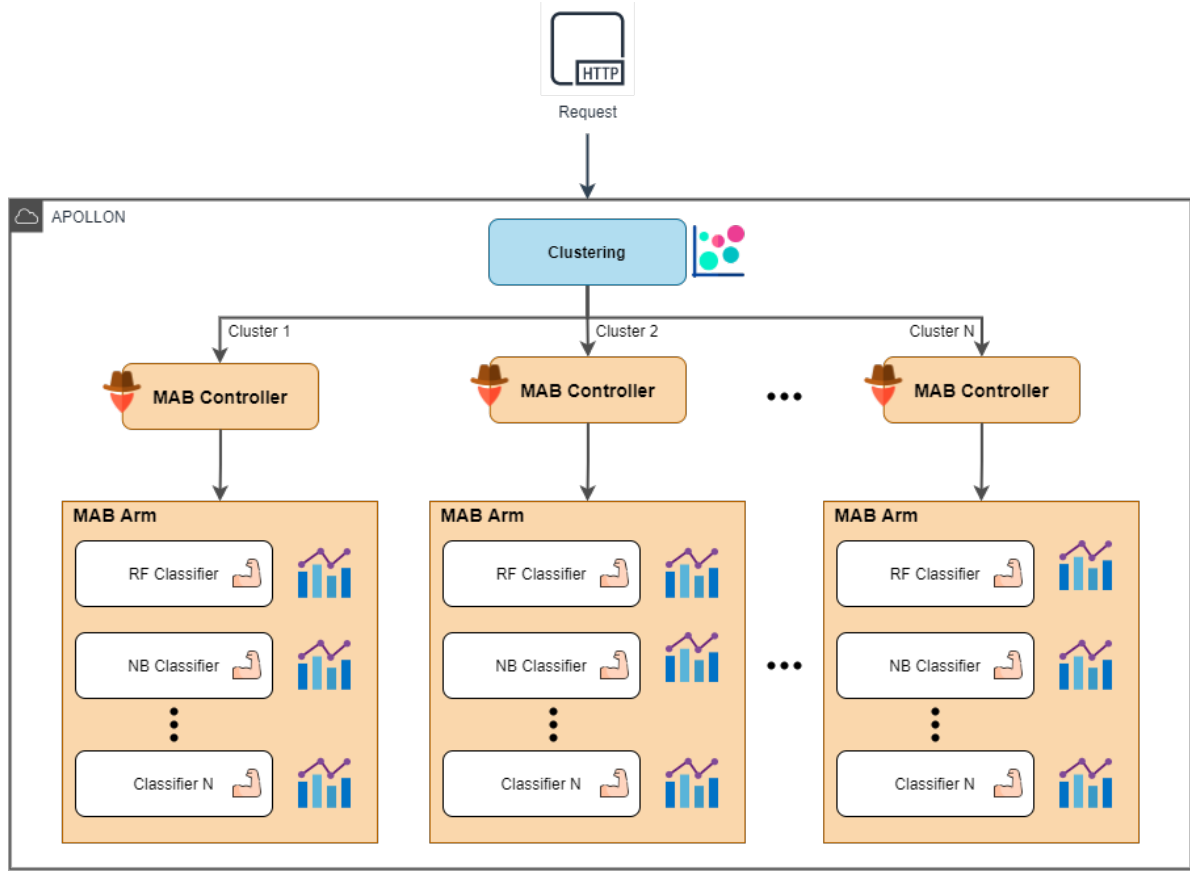
Adversarial training has been extensively studied in the field of visual computing. Goodfellow et al. [58] showed that by re-training a neural network on a dataset that includes both original and adversarial samples, the network's effectiveness in countering adversarial samples can be significantly enhanced.

However, despite its potential benefits, adversarial training is not without limitations. Over the years, researchers have made efforts to address these shortcomings and enhance its effectiveness.

One notable challenge with adversarial training is its computational and time requirements. The process can be resource-intensive, making it impractical for systems with limited capabilities or real-time applications. Additionally, while adversarial training can bolster a model's resilience against certain attacks, it may not provide comprehensive defense against all types of adversarial techniques.

Obtaining a sufficient number of high-quality adversarial samples for training purposes can also be a hurdle. Generating diverse and representative adversarial samples can be challenging, impacting the efficacy of the training process [59]. Moreover, the scarcity of such samples can hinder the scalability and generalizability of adversarial training.

Another consideration is the potential trade-off between adversarial robustness and accuracy on clean data [59]. Adversarial training has been observed to cause a decrease in a model's performance on non-adversarial inputs. This trade-off poses a dilemma, as maintaining high accuracy on clean data is crucial in many real-world scenarios.


 Figure 2: *Apollon* Architecture

3.4.3. Adversarial detection defenses

An alternative approach to combat adversarial attacks involves the implementation of detection mechanisms capable of identifying the presence of adversarial samples [60]. These mechanisms utilize various techniques such as direct classification, neural network uncertainty, or input processing.

However, it has been observed that these detection mechanisms are often insufficient in effectively defending against AML attacks [60]. Despite their potential, they exhibit weaknesses that can be exploited by sophisticated adversarial techniques [56].

The weaknesses of these detection mechanisms stem from the constantly evolving nature of adversarial attacks. Adversarial samples are intentionally crafted to deceive the detection mechanisms, making it challenging to accurately distinguish between genuine and adversarial inputs.

4. Apollon

In this paper, we propose a robust defence system called *Apollon*, which is designed to protect an IDS against AML attacks. *Apollon* is composed of multiple layers to provide better security than traditional IDS and previous works that rely solely on training with adversarial traffic. The proposed system combines multiple classifiers, a *Multi-Armed Bandits*

(*MAB*) algorithm, and requests clustering to provide robust defence against AML attacks.

The first layer of *Apollon* involves using multiple classifiers instead of a single classifier that is traditionally used in IDS. The concept behind utilizing multiple classifiers is to increase the difficulty for potential attackers attempting to replicate the IDS model. This is because they cannot predict which specific model will be responsible for classifying a given request.

To dynamically select the optimal classifier or set of classifiers for each input, *Apollon* involves using a *Multi-Armed Bandits (MAB)* with *Thompson sampling*. The *MAB* is responsible for selecting the arm (classifier) to use for each request based on the current state of the system.

Finally, requests are clustered, and there is a version of each classifier for each cluster, trained only with the information of that cluster. Clustering is used to add another layer of uncertainty to the system, as the attacker cannot predict in a simple way which cluster a request belongs to.

With *Apollon*, we are able to maintain the performance of traditional IDSs when it comes to network traffic without AML attacks. We achieve this by utilizing the best classifiers for each type of request. Additionally, *Apollon* offers a solution to prevent attackers from easily learning from the behavior of our classifiers through AML techniques.

Figure 2 shows the architecture of *Apollon*, and the flow that a network traffic request follows until it is classified.

Listing 1 shows the output of an *Apollon* training process example with two clusters and three classifiers. As can be seen, the training data is first divided between the two clusters, the classifiers are trained for each cluster and finally, the rewards are assigned. The Listing 2 shows a selection process example of the classifier for each request, the predicted value and the real value.

```
Cluster 0: Len of request 274888
Training arm 0 on cluster 0
Training arm 1 on cluster 0
Training arm 2 on cluster 0
Cluster 1: Len of request 358525
Training arm 0 on cluster 1
Training arm 1 on cluster 1
Training arm 2 on cluster 1

Setting reward_sums arm 0 on cluster 0
Setting reward_sums arm 1 on cluster 0
Setting reward_sums arm 2 on cluster 0
Setting reward_sums arm 0 on cluster 1
Setting reward_sums arm 1 on cluster 1
Setting reward_sums arm 2 on cluster 1
```

Listing 1: *Apollon* training process example

Selected arm: 0.0	Predicted:0	Actual:0
Selected arm: 0.0	Predicted:0	Actual:0
Selected arm: 2.0	Predicted:0	Actual:0
Selected arm: 0.0	Predicted:0	Actual:0
Selected arm: 1.0	Predicted:1	Actual:1
Selected arm: 2.0	Predicted:0	Actual:0

Listing 2: *Apollon* classifier selection process example

In the following, each of the layers that influence the final classification of a network traffic request will be detailed.

4.1. Multiple Classifiers

The first layer of *Apollon* involves using multiple classifiers instead of a single classifier, which is typically used in traditional IDS. The idea of using multiple classifiers is to make it more difficult for the attacker to replicate the IDS model, as he is not able to predict which model is going to classify the request. By utilizing a greater number of diverse classifiers, the system becomes more resilient by introducing greater uncertainty, and better classifiers implies better performance in the *Apollon* system score.

In *Apollon*, we can use any type of classifier. These classifiers can be either the most common ones based on Deep Learning or Machine Learning, or classifiers based on network traffic request forecasting techniques, or the more classical ones based on rule systems.

4.2. Multi-Armed Bandit

The second layer of the *Apollon* defence system involves the use of a *Multi-Armed Bandits* (MAB) algorithm to select the appropriate classifier or set of classifiers for each network traffic request. The MAB is responsible for selecting the best

classifier or set of classifiers to evaluate whether a request is benign or malicious. This approach avoids the need for manual tuning of thresholds or weights for each classifier.

The MAB algorithm works by selecting the arm, or classifier, that has the highest probability of providing the correct classification. In *Apollon*, we use *Thomson Sampling*, which is a popular algorithm for solving the MAB problem. Thomson Sampling balances exploration and exploitation of the available classifiers, ensuring that the system selects the optimal classifier or set of classifiers while still being responsive to new and unknown types of traffic.

The MAB algorithm in *Apollon* is designed to take into account the different types of classifiers used in the system. For instance, if the *Random Forest* classifier has a high probability of being correct, but the *Naive Bayes* and *Logistic Regression* classifiers have lower probabilities, the MAB algorithm will select the *Random Forest* classifier for that particular request. In this way, the MAB algorithm ensures that the system selects the optimal set of classifiers for each request, improving the overall accuracy of the classification.

The MAB algorithm is constantly updating the probabilities of the different classifiers based on their previous performance. Thus, the system can adapt to changes in the traffic patterns over time, ensuring that the system is always up to date with the latest types of attacks. This update process is described in Algorithm 1. Here, S_i and F_i are the number of observed successes and failures for arm i , respectively, and θ_{θ_i} is the estimated probability of obtaining a positive reward from arm i . The algorithm uses the Beta distribution as an a priori distribution for the parameters θ_{θ_i} , and updates it with the observed data using Bayes' theorem. The algorithm chooses the arm (ML classifier) that has the highest probability of being the best according to the samples from the posterior distribution.

Algorithm 1 *Apollon* Thompson Sampling

```
Init  $S_i = 0$  y  $F_i = 0$  for each arm  $i$ 
for  $t = 1, 2, \dots$  do
    For each arm  $i$ , sample  $\theta_i$  of Beta distribution  $(S_i + 1, F_i + 1)$ 
    Choose the arm  $I_t$  that maximizes  $\theta_i$ 
    Observe the reward  $X_t$  of the arm  $I_t$ .
    if  $X_t = 1$  then
        Increment  $S_{I_t}$  by one
    else
        Increment  $F_{I_t}$  by one
    end if
end for
```

By using a *Multi-Armed Bandits* algorithm, *Apollon* can dynamically select the optimal classifier or set of classifiers for each network traffic request, making the system more responsive to new types of attacks. The use of *Thomson Sampling* ensures that the system is balanced between exploration and exploitation, improving the overall attacks detection rate of the classification. Figure 3 shows the diagram of the MAB algorithm with multiple classifiers.

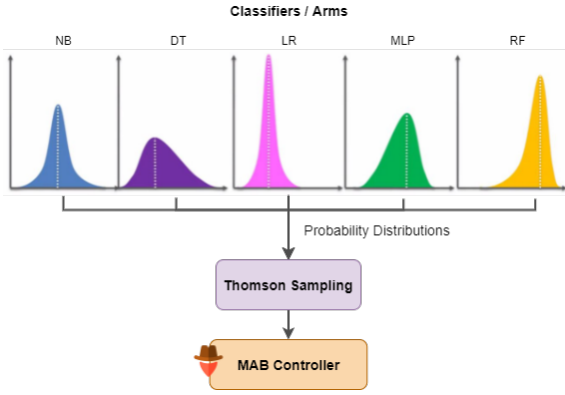


Figure 3: Apollon *Multi-Armed Bandits* algorithm with multiple classifiers

4.3. Traffic requests clustering

The final layer of the *Apollon* defence system involves clustering the network traffic requests based on their features, and then training a separate version of each classifier for each cluster. By doing this, the system is able to achieve higher accuracy for each cluster by training the classifiers specifically for the traffic patterns in that cluster.

The traffic requests are clustered using the *K-Means* algorithm [61], which is a popular clustering algorithm that is used in many Machine Learning applications. The *K-Means* algorithm works by randomly selecting k points as the initial centroids, and then iteratively updating the centroids until the clusters converge. In *Apollon*, we use the *K-Means* algorithm to cluster the network traffic requests based on their features, ensuring that requests with similar features are grouped together. These features are the same ones used by the classifiers to evaluate the requests, ensuring that the clustering is based on the same information that the classifiers use to make their decisions.

For each cluster, a separate version of each classifier is trained using only the traffic requests in that cluster. This ensures that each classifier is specifically tuned to the traffic patterns in that cluster.

When a new network traffic request arrives at *Apollon*, it is classified into the appropriate cluster based on its features. The *Multi-Armed Bandits* algorithm is then used to select the optimal set of classifiers for that cluster, taking into account the performance of each classifier in that specific cluster. Then, the selected classifier or set of classifiers evaluate the request and determine whether it is benign or malicious.

By clustering the network traffic requests and training separate versions of each classifier for each cluster, *Apollon* allows the *Multi-Armed Bandit* algorithm to generate multiple probability distributions for each classifier, depending on the type of request received. This combination of techniques makes it challenging for potential attackers to identify the classifier providing the response, reducing the likelihood of successful imitation.

5. Evaluation

In this section, we present the methodology and results of our evaluation, which assesses the performance of *Apollon* in traditional and AML attack environments. The code that was used and created during the evaluation of our proposal is completely open and accessible. It can be found on GitHub¹.

5.1. Methodology

We designed two scenarios to assess the performance of our proposed solution. In the first scenario, we tested our solution on three datasets: CIC-IDS-2017, CSE-CIC-IDS-2018, and CIC-DDoS-2019. Due to our lack of powerful machines, we opted to use a representative subset of the datasets instead of the complete data to efficiently train and test our models. For the dataset CIC-DDoS-2019 the whole dataset has been used, while for the dataset CSE-CIC-IDS-2018 the subset from 02-15-2018 has been used. Finally, for the CIC-IDS-2017 dataset, the following subsets of data have been selected:

- *Friday WorkingHours Afternoon DDoS*
- *Friday WorkingHours Afternoon PortScan*
- *Friday WorkingHours Morning*
- *Monday WorkingHours*
- *Thursday WorkingHours Afternoon Infiltration*
- *Thursday WorkingHours Morning WebAttacks*
- *Tuesday WorkingHours*

We compared the results of our solution with the classifiers extracted from related work. In the second scenario, we used several grey/black-box Adversarial Machine Learning attacks to evaluate the ability of our solution to defend against such attacks.

The classifiers scores used to compare our solution may be different than those reported in related work due to several factors. Firstly, the machine on which the training is performed may differ from that of related work. This can impact the speed and efficiency of the training process, which in turn can affect the final accuracy of the classifiers. Secondly, the pre-processing of the data may be different between our solution and related work. Pre-processing techniques can greatly impact the quality of the data and hence the performance of the classifiers. Therefore, differences in pre-processing techniques can lead to varying levels of accuracy in the classifiers. It is important to take into account these differences when comparing our solution to related work, and to consider the impact of these factors on the performance of the classifiers.

All the experiments were performed on a *Ubuntu 20.04.5 LTS* machine with an *Intel(R) Core(TM) i7-7700 CPU @ 3.60GHz* processor and 16 GB of RAM memory.

¹<https://github.com/antonioalfa22/apollon>

5.1.1. Traditional Network Traffic and Attacks

In this test scenario, we utilized the CIC-IDS-2017, CSE-CIC-IDS-2018, and CIC-DDoS-2019 datasets to train a set of classifiers to compare with our proposed solution. Our goal was to assess the performance of our solution with the traditional network traffic and web attacks.

We have utilized default hyperparameters of the classifiers because the objective of this scenario is not to maximize the performance of these classifiers but rather to compare them with our solution. The classifiers used in this scenario are the following:

- *Multilayer Perceptron (MLP)*: hidden_layer_sizes = (32), max_iter = 200
- *Random Forest (RF)*: n_estimators = 100
- *Decision Trees (DT)*
- *Naive Bayes (NB)*
- *Logistic Regression (LR)*

5.1.2. Adversarial Machine Learning Attacks

In this test scenario, we used several grey/black-box Adversarial Machine Learning attacks to evaluate the ability of our solution to defend against such attacks. To simplify the evaluation process, we used the CIC-IDS-2017 dataset to train the classifiers and our proposed solution because it is the most popular dataset. We compare the accuracy and the detection rate of the classifiers and our proposed solution against the grey/black-box Adversarial Machine Learning attacks. As we mentioned before, we only test with grey/black-box attacks because white-box attacks are not realistic in real-world scenarios.

The attacks used in this scenario are the following:

- *Zeroth-order optimization attack (ZOO)* [50]
- *HopSkipJump attack (HSJA)* [52]
- *W-GAN based attacks* [6]

We opted for these particular attacks because they encompass a broad spectrum of potential Adversarial Machine Learning evasion strategies and are among the most widely used and successful.

5.2. Results

Throughout the development of the evaluation we have decided to set a seed, so that the results can be replicated: the *seed* = 42.

Before training the classifiers, a common pre-processing step was performed on the data from all datasets. This step is essential in standardizing the datasets, ensuring that we are working uniformly with each of them. To achieve this, a combination of *sklearn* functions such as *RobustScaler* and *Normalizer* was utilized. To make features less sensitive to outliers, the *RobustScaler* subtracts the median and adjusts the data based on the quantile range. The *Normalizer* scales

the input data set to have a norm of 1 and values between 0 and 1.

In addition to standardizing the datasets, additional steps were taken to further prepare the data for training our classifiers. We avoided excessive pre-processing, to preserve as much data as possible, so we only eliminated duplicate elements and attributes that had only unique data.

We also encountered some data with missing values (*NaN*) and explored various methods to handle this issue: filling in missing values with a fixed value, with the mean or median of the feature's non-missing values, with a prediction based on the other features and with a prediction based on *k*-nearest neighbors. Our experiments revealed that using a constant value of 0 to fill in missing values achieved the best performance in the training step.

We apply this pre-processing to all datasets equally. When validating the results we used a cross validation, in which the following parameters have been used:

- n_splits: 5
- test_size: 0.2
- random_state: seed

By using cross-validation we ensure that the results obtained are representative and not the result of an isolated iteration.

The results of our experiments in the two evaluation scenarios are presented below.

5.2.1. Traditional Network Traffic and Attacks

In this scenario, we trained a the selected classifiers on the CIC-IDS-2017, CSE-CIC-IDS-2018, and CIC-DDoS-2019 datasets. Additionally, we trained our proposed solution on the same datasets and with the same classifiers. To train *Apollon*, we use the K-Means [62] algorithm to cluster the data into 2 clusters. For each cluster, we train the selected classifiers and we update the *Multi-Armed Bandits* algorithm with the results.

The results obtained from the experiments are presented in the tables 2, 3, and 4. Based on the tables, our solution demonstrates high detection rate and accuracy scores, comparable to the classifiers chosen for comparison. It's important to mention that among all the datasets, *Apollon* doesn't achieve the best or the worst scores. This is due to the fact that *Apollon* internally selects from the same classifiers. Hence, the highest score that *Apollon* can achieve is limited to the best classifier's maximum score, while it can never perform as poorly as the worst classifier since it has other better options.

Our results demonstrate that even with the integration of new security mechanisms, *Apollon* can still provide high accuracy and detection rate scores in traditional network traffic classification environments. Therefore, *Apollon* is capable of preserving the fundamental functionality of an IDS.

5.2.2. Adversarial Machine Learning Attacks

In this scenario, we have launched three types of Adversarial Machine Learning attacks on the selected classifiers

Metrics	CIC-IDS-2017					
	MLP	NB	RF	DT	LR	Apollon
Accuracy	0.9766	0.6694	0.9996	0.9980	0.9516	0.9740
Detection rate	0.9731	0.8049	0.9996	0.9962	0.8360	0.9420
F1	0.9760	0.6118	0.9991	0.9959	0.8858	0.7699
ROC	0.9998	0.8289	1.0000	0.9972	0.9902	0.9420

Table 2

Results of the traditional network traffic and attacks scenario on the CIC-IDS-2017 dataset.

Metrics	CSE-CIC-IDS-2018					
	MLP	NB	RF	DT	LR	Apollon
Accuracy	0.9916	0.7280	0.9993	0.9939	0.9593	0.9064
Detection rate	0.9367	0.8563	0.9967	0.9961	0.6008	0.9419
F1	0.9545	0.5507	0.9963	0.9968	0.6556	0.7303
ROC	0.9945	0.8605	0.9993	0.9984	0.9592	0.9419

Table 3

Results of the traditional network traffic and attacks scenario on the CSE-CIC-IDS-2018 dataset.

Metrics	CIC-DDoS-2019					
	MLP	NB	RF	DT	LR	Apollon
Accuracy	0.9998	0.9990	0.9999	0.9999	0.9970	0.9996
Detection rate	0.8985	0.8709	0.9932	0.9999	0.7930	0.9691
F1	0.9186	0.7148	0.9957	0.9991	0.8541	0.8521
ROC	0.9817	0.9753	0.9999	0.9999	0.9796	0.9691

Table 4

Results of the traditional network traffic and attacks scenario on the CIC-DDoS-2019 dataset.

and against our solution. These attacks are *Zeroth-order optimization attack (ZOO)*, *HopSkipJump attack (HSJA)* and *W-GAN* based attacks.

The classifiers and the *Apollon* implementation used as targets of the attacks are the ones trained in the previous environment with the CIC-IDS-2017 dataset.

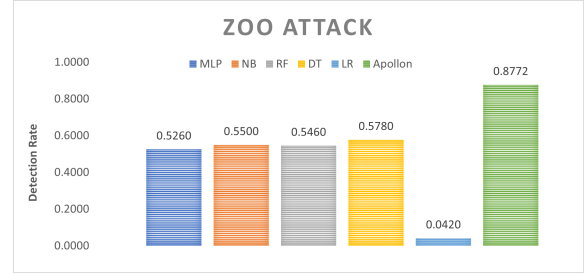
Starting with the *Zeroth-order optimization attack (ZOO)*, we have used the open source implementation provided by ART [11], and created a Classifier class so that *Apollon* can be used as a model. The attack was launched with the following parameters for each classifier:

- **classifier**: the Classifier class instance with the classifier to be attacked.
- **targeted**: True.
- **learning_rate**: 0.01.

Metrics	ZOO attack					
	MLP	NB	RF	DT	LR	Apollon
Accuracy	0.7630	0.7160	0.7730	0.7890	0.5210	0.9304
Detection rate	0.5260	0.5500	0.5460	0.5780	0.0420	0.8772

Table 5

Results of the ZOO AML attack


Figure 4: ZOO attack detection rate results

- **max_iter**: 100.

The attack was launched against the classifiers and the results are shown in Table 5 and in the Figure 4. According to the findings, the attack was effective in every scenario, resulting in reduced detection rates across all classifiers. Nonetheless, even though the *Apollon* implementation's accuracy and detection rate scores also declined, they remained considerably superior to those of the other classifiers, with high accuracy and detection rate scores.

To launch the *HopSkipJump attack (HSJA)*, we have used the open source implementation provided by ART, and created a Classifier as in the ZOO attack. The attack was launched with the following parameters for each classifier:

- **classifier**: the Classifier class instance with the classifier to be attacked.
- **targeted**: True.
- **max_iter**: 100.
- **norm**: inf.

The results of the attack are shown in Table 6 and in the Figure 5. The attack was very effective in all the classifiers, resulting in reduced detection rates to scores close to zero. The exception was the *Apollon* implementation, which was able to maintain a detection rate > 0.5.

Finally, the *Wasserstein Generative Adversarial Network (W-GAN)* based attack was launched with a custom implementation available in the *Apollon* repository on GitHub. This implementation was based on the *IDSGAN* attack [6], updated to the needs of the selected dataset. The results in Table 7 and in the Figure 6 were obtained after launching the attack with 100 epochs.

The W-GAN based attack was the most effective attack, reducing the detection rate to zero in all the classifiers. Our

Metrics	HopSkipJump attack					
	MLP	NB	RF	DT	LR	Apollon
Accuracy	0.5002	0.4301	0.5001	0.5051	0.5900	0.7550
Detection Rate	0.0000	0.0000	0.0000	0.0100	0.1800	0.5260

Table 6
Results of the *HopSkipJump* AML attack

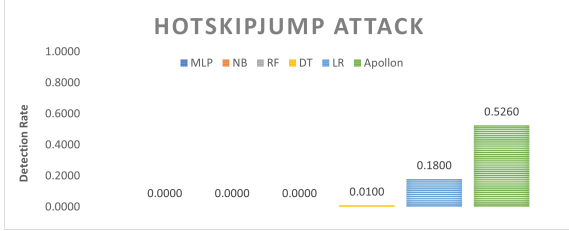


Figure 5: *HopSkipJump* attack Detection Rate results

Metrics	W-GAN based attack					
	MLP	NB	RF	DT	LR	Apollon
Accuracy	0.5002	0.4398	0.5001	0.4280	0.3230	0.6260
Detection Rate	0.0000	0.0000	0.0000	0.0000	0.0000	0.4250

Table 7
Results of the *W-GAN* based AML attack

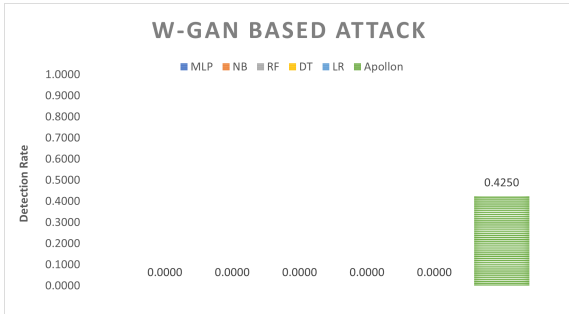


Figure 6: *W-GAN* based attack detection rate results

solution, on the other hand, was able to maintain a detection rate > 0.4 .

Apollon's improved accuracy and detection rates in AML attacks can be attributed to the inclusion of an uncertainty component in its model selection process. This renders it difficult to train a model solely based on its responses.

The experimental results of the scenario reveals that our solution exhibits greater robustness in comparison to the other classifiers used individually.

However, we also observed that while our solution effectively reduces the effectiveness of the attacks, it does not completely nullify them. This means that there is still room for improvement in terms of enhancing the solution's robustness to further strengthen its resistance against such attacks.

In particular, if we were to generate the attacks with more time, such as by increasing the number of iterations or epochs, it is likely that the effectiveness of these attacks against our solution would increase.

6. Conclusions and Future Work

In conclusion, this paper presents *Apollon*, a new robust defense system against Adversarial Machine Learning attacks on Intrusion Detection Systems. *Apollon* utilizes a *Multi-Armed Bandits* model to select the best-suited classifier in real-time for each input with Thompson sampling, adding a layer of uncertainty to the IDS behavior, that makes it more difficult for attackers to replicate the IDS and generate adversarial traffic.

Our experimental evaluation on several datasets shows that *Apollon* can successfully detect attacks without compromising its performance on normal network traffic data, and can prevent attackers from learning the IDS behavior in realistic training times. These results demonstrate that *Apollon* is an effective defense system against AML attacks in IDS, which can help to enhance the security of critical systems. Nevertheless, *Apollon* does not completely eliminate the risk of AML attacks, only mitigates it increasing the time and effort required by attackers to generate adversarial traffic.

With the aim of improving the performance and robustness of *Apollon*, we plan to explore the use of other *MAB* models and implementations, such as Bayesian Optimization or Deep Bayesian Bandits. We also plan to explore the use of other classifiers and ML models, such as requests forecasting models, which can be used to predict the expected number of requests in the next time window and compare it with the actual number of requests. Finally, we plan to explore the use of other datasets, to evaluate the performance of *Apollon* in different network environments.

References

- [1] A. Thakkar and R. Lohiya, "A review of the advancement in intrusion detection datasets," *Procedia Computer Science*, vol. 167, pp. 636–645, 2020.
- [2] E. E. Abdallah, A. F. Otoom, *et al.*, "Intrusion detection systems using supervised machine learning techniques: A survey," *Procedia Computer Science*, vol. 201, pp. 205–212, 2022.
- [3] Z. K. Maseer, R. Yusof, N. Bahaman, S. A. Mostafa, and C. F. M. Foozy, "Benchmarking of machine learning for anomaly based intrusion detection systems in the cids2017 dataset," *IEEE access*, vol. 9, pp. 22351–22370, 2021.
- [4] L. Huang, A. D. Joseph, B. Nelson, B. I. Rubinstein, and J. D. Tygar, "Adversarial machine learning," in *Proceedings of the 4th ACM workshop on Security and artificial intelligence*, pp. 43–58, 2011.
- [5] S. Zhao, J. Li, J. Wang, Z. Zhang, L. Zhu, and Y. Zhang, "attackgan: Adversarial attack against black-box ids using generative adversarial networks," *Procedia Computer Science*, vol. 187, pp. 128–133, 2021.
- [6] Z. Lin, Y. Shi, and Z. Xue, "Idsgan: Generative adversarial networks for attack generation against intrusion detection," in *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pp. 79–91, Springer, 2022.
- [7] G. Liu, W. Zhang, X. Li, K. Fan, and S. Yu, "Vulnergan: a backdoor attack through vulnerability amplification against machine learning-based network intrusion detection systems," *Science China Information Sciences*, vol. 65, no. 7, pp. 1–19, 2022.

- [8] P. T. Duy, N. H. Khoa, A. G.-T. Nguyen, V.-H. Pham, *et al.*, “Digfupas: Deceive ids with gan and function-preserving on adversarial samples in sdn-enabled networks,” *Computers & Security*, vol. 109, p. 102367, 2021.
- [9] B. Mukherjee, L. T. Heberlein, and K. N. Levitt, “Network intrusion detection,” *IEEE network*, vol. 8, no. 3, pp. 26–41, 1994.
- [10] A. Pharate, H. Bhat, V. Shilimkar, and N. Mhetre, “Classification of intrusion detection system,” *International Journal of Computer Applications*, vol. 118, no. 7, 2015.
- [11] M.-I. Nicolae, M. Sinn, M. N. Tran, B. Buesser, A. Rawat, M. Wistuba, V. Zantedeschi, N. Baracaldo, B. Chen, H. Ludwig, I. Molloy, and B. Edwards, “Adversarial robustness toolbox v1.2.0,” *CoRR*, vol. 1807.01069, 2018.
- [12] E. De Cristofaro, “An overview of privacy in machine learning,” *arXiv preprint arXiv:2005.08679*, 2020.
- [13] V. Kuleshov and D. Precup, “Algorithms for multi-armed bandit problems,” *arXiv preprint arXiv:1402.6028*, 2014.
- [14] M. C. Machado, S. Srinivasan, and M. Bowling, “Domain-independent optimistic initialization for reinforcement learning,” *arXiv preprint arXiv:1410.4604*, 2014.
- [15] A. Carpentier, A. Lazaric, M. Ghavamzadeh, R. Munos, and P. Auer, “Upper-confidence-bound algorithms for active learning in multi-armed bandits,” in *Algorithmic Learning Theory: 22nd International Conference, ALT 2011, Espoo, Finland, October 5-7, 2011. Proceedings* 22, pp. 189–203, Springer, 2011.
- [16] S. Agrawal and N. Goyal, “Analysis of thompson sampling for the multi-armed bandit problem,” in *Conference on learning theory*, pp. 39–1, JMLR Workshop and Conference Proceedings, 2012.
- [17] H. Park and M. K. S. Faradonbeh, “Analysis of thompson sampling for partially observable contextual multi-armed bandits,” *IEEE Control Systems Letters*, vol. 6, pp. 2150–2155, 2021.
- [18] D. Ramos, P. Faria, L. Gomes, P. Campos, and Z. Vale, “Selection of features in reinforcement learning applied to energy consumption forecast in buildings according to different contexts,” *Energy Reports*, vol. 8, pp. 423–429, 2022.
- [19] D. Ramos, P. Faria, L. Gomes, P. Campos, and Z. Vale, “A learning approach to improve the selection of forecasting algorithms in an office building in different contexts,” in *Progress in Artificial Intelligence: 21st EPIA Conference on Artificial Intelligence, EPIA 2022, Lisbon, Portugal, August 31–September 2, 2022, Proceedings*, pp. 271–281, Springer, 2022.
- [20] M. V. Mahoney and P. K. Chan, “An analysis of the 1999 darpa/lincoln laboratory evaluation data for network anomaly detection,” in *International Workshop on Recent Advances in Intrusion Detection*, pp. 220–237, Springer, 2003.
- [21] M. Tavallaei, E. Bagheri, W. Lu, and A. A. Ghorbani, “A detailed analysis of the kdd cup 99 data set,” in *2009 IEEE Symposium on Computational Intelligence for Security and Defense Applications*, pp. 1–6, 2009.
- [22] J. M. Hugh, “Testing intrusion detection systems: A critique of the 1998 and 1999 darpa intrusion detection system evaluations as performed by lincoln laboratory,” *ACM Transactions on Information and System Security*, 2000.
- [23] A. M. A. Tobi and I. Duncan, “Kdd 1999 generation faults: a review and analysis,” *Journal of Cyber Security Technology*, vol. 2, no. 3–4, pp. 164–200, 2018.
- [24] I. Sharafaldin, A. H. Lashkari, and A. A. Ghorbani, “Toward generating a new intrusion detection dataset and intrusion traffic characterization,” *ICISSp*, vol. 1, pp. 108–116, 2018.
- [25] M. Ring, S. Wunderlich, D. Scheuring, D. Landes, and A. Hotho, “A survey of network-based intrusion detection data sets,” *Computers & Security*, vol. 86, pp. 147–167, 2019.
- [26] J. Shroff, R. Walambe, S. K. Singh, and K. Kotecha, “Enhanced security against volumetric ddos attacks using adversarial machine learning,” *Wireless Communications and Mobile Computing*, vol. 2022, 2022.
- [27] I. Sharafaldin, A. H. Lashkari, and A. A. Ghorbani, “Toward generating a new intrusion detection dataset and intrusion traffic characterization,” in *International Conference on Information Systems Security and Privacy*, 2018.
- [28] M. Pujari, Y. Pacheco, B. Cherukuri, and W. Sun, “A comparative study on the impact of adversarial machine learning attacks on contemporary intrusion detection datasets,” *SN Computer Science*, vol. 3, no. 5, pp. 1–12, 2022.
- [29] I. Sharafaldin, A. H. Lashkari, S. Hakak, and A. A. Ghorbani, “Developing realistic distributed denial of service (ddos) attack dataset and taxonomy,” in *2019 International Carnahan Conference on Security Technology (ICCST)*, pp. 1–8, 2019.
- [30] B. Thiyyam and S. Dey, “Efficient feature evaluation approach for a class-imbalanced dataset using machine learning,” *Procedia Computer Science*, vol. 218, pp. 2520–2532, 2023. International Conference on Machine Learning and Data Engineering.
- [31] M. Akshay Kumar, D. Samiyya, P. M. D. R. Vincent, K. Srinivasan, C.-Y. Chang, and H. Ganesh, “A hybrid framework for intrusion detection in healthcare systems using deep learning,” *Frontiers in Public Health*, vol. 9, 2022.
- [32] S. Huang and K. Lei, “Igan-ids: An imbalanced generative adversarial network towards intrusion detection system in ad-hoc networks,” *Ad Hoc Networks*, vol. 105, p. 102177, 2020.
- [33] Z. Wu, H. Zhang, P. Wang, and Z. Sun, “Rtids: a robust transformer-based approach for intrusion detection system,” *IEEE Access*, vol. 10, pp. 64375–64387, 2022.
- [34] R. Abdulhammed, M. Faezipour, H. Musaffer, and A. Abuzneid, “Efficient network intrusion detection using pca-based dimensionality reduction of features,” in *2019 International Symposium on Networks, Computers and Communications (ISNCC)*, pp. 1–6, 2019.
- [35] O. Faker and E. Dogdu, “Intrusion detection using big data and deep learning techniques,” in *Proceedings of the 2019 ACM Southeast Conference*, ACM SE ’19, (New York, NY, USA), p. 86–93, Association for Computing Machinery, 2019.
- [36] R. E. Wright, “Logistic regression,” 1995.
- [37] R. A. R. Ashfaq, X.-Z. Wang, J. Z. Huang, H. Abbas, and Y.-L. He, “Fuzziness based semi-supervised learning approach for intrusion detection system,” *Information sciences*, vol. 378, pp. 484–497, 2017.
- [38] A. Cutler, D. R. Cutler, and J. R. Stevens, “Random forests,” *Ensemble machine learning: Methods and applications*, pp. 157–175, 2012.
- [39] L. Rokach and O. Maimon, “Decision trees,” *Data mining and knowledge discovery handbook*, pp. 165–192, 2005.
- [40] S. Suthaharan and S. Suthaharan, “Support vector machine,” *Machine learning models and algorithms for big data classification: thinking with examples for effective learning*, pp. 207–235, 2016.
- [41] J. Zhang, M. Zulkernine, and A. Haque, “Random-forests-based network intrusion detection systems,” *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 38, no. 5, pp. 649–659, 2008.
- [42] N. B. Amor, S. Benferhat, and Z. Elouedi, “Naive bayes vs decision trees in intrusion detection systems,” in *Proceedings of the 2004 ACM symposium on Applied computing*, pp. 420–424, 2004.
- [43] I. J. Goodfellow, J. Shlens, and C. Szegedy, “Explaining and harnessing adversarial examples,” *arXiv preprint arXiv:1412.6572*, 2014.
- [44] S.-M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard, “Deepfool: a simple and accurate method to fool deep neural networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2574–2582, 2016.
- [45] N. Carlini and D. Wagner, “Towards evaluating the robustness of neural networks,” in *2017 IEEE Symposium on Security and Privacy (SP)*, pp. 39–57, IEEE, 2017.
- [46] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, and A. Swami, “The limitations of deep learning in adversarial settings,” in *2016 IEEE European symposium on security and privacy (EuroS&P)*, pp. 372–387, IEEE, 2016.
- [47] A. Kurakin, I. Goodfellow, and S. Bengio, “Adversarial machine learning at scale,” *arXiv preprint arXiv:1611.01236*, 2016.

- [48] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," *arXiv preprint arXiv:1706.06083*, 2017.
- [49] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," *Communications of the ACM*, vol. 63, no. 11, pp. 139–144, 2020.
- [50] P.-Y. Chen, H. Zhang, Y. Sharma, J. Yi, and C.-J. Hsieh, "Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models," in *Proceedings of the 10th ACM workshop on artificial intelligence and security*, pp. 15–26, 2017.
- [51] J. Chen and M. I. Jordan, "Boundary attack++: Query-efficient decision-based adversarial attack," *arXiv preprint arXiv:1904.02144*, vol. 2, no. 7, 2019.
- [52] J. Chen, M. I. Jordan, and M. J. Wainwright, "Hopskipjumpattack: A query-efficient decision-based attack," 2020.
- [53] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville, "Improved training of wasserstein gans," *Advances in neural information processing systems*, vol. 30, 2017.
- [54] A. Alotaibi and M. A. Rassam, "Adversarial machine learning attacks against intrusion detection systems: A survey on strategies and defense," *Future Internet*, vol. 15, no. 2, p. 62, 2023.
- [55] C. Kou, H. K. Lee, E.-C. Chang, and T. K. Ng, "Enhancing transformation-based defenses using a distribution classifier," *arXiv preprint arXiv:1906.00258*, 2019.
- [56] A. Athalye, N. Carlini, and D. Wagner, "Obfuscated gradients give a false sense of security, circumventing defenses to adversarial examples," in *International conference on machine learning*, pp. 274–283, PMLR, 2018.
- [57] H. Xu, Y. Ma, H.-C. Liu, D. Deb, H. Liu, J.-L. Tang, and A. K. Jain, "Adversarial attacks and defenses in images, graphs and text: A review," *International Journal of Automation and Computing*, vol. 17, pp. 151–178, 2020.
- [58] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *arXiv preprint arXiv:1412.6572*, 2014.
- [59] T. Bai, J. Luo, J. Zhao, B. Wen, and Q. Wang, "Recent advances in adversarial training for adversarial robustness," *arXiv preprint arXiv:2102.01356*, 2021.
- [60] G. Zizzo, C. Hankin, S. Maffei, and K. Jones, "Adversarial machine learning beyond the image domain," in *Proceedings of the 56th Annual Design Automation Conference 2019*, pp. 1–4, 2019.
- [61] K. P. Sinaga and M.-S. Yang, "Unsupervised k-means clustering algorithm," *IEEE access*, vol. 8, pp. 80716–80727, 2020.
- [62] A. Likas, N. Vlassis, and J. J. Verbeek, "The global k-means clustering algorithm," *Pattern recognition*, vol. 36, no. 2, pp. 451–461, 2003.