

RESEARCH ARTICLE

A fast and accurate threat detection and prevention architecture using stream processing

Antonio G. Pastana Lobato¹  | Martin Andreoni Lopez^{1,2}  | Alvaro A. Cardenas³  |
Otto Carlos M. B. Duarte¹  | Guy Pujolle² 

¹GTA/COPPE/UFRJ, Universidade Federal do Rio de Janeiro, Rio de Janeiro, Brazil

²Laboratoire d'Informatique de Paris 6, CNRS, Sorbonne Université, Paris, France

³Department of Computer Science and Engineering, University of California, Santa Cruz, Santa Cruz, CA, USA

Correspondence

Otto Carlos M. B. Duarte, GTA/COPPE/UFRJ, Universidade Federal do Rio de Janeiro, Rio de Janeiro, Brazil.

Email: otto@gta.ufrj.br

Summary

Late detection of security breaches increases the risk of irreparable damages and limits any mitigation attempts. We propose a fast and accurate threat detection and prevention architecture that combines the advantages of real-time streaming with batch processing over a historical database. We create a dataset by capturing both legitimate and malicious traffic and propose two ways of combining packets into flows, one considering a time window and the other analyzing the first few packets of each flow per period. We also investigate the effectiveness of our proposal on real-world network traces obtained from a significant Brazilian network operator providing broadband Internet to their customers. We implement and evaluate three classification algorithms and two anomaly detection methods. The results show an accuracy higher than 95% and an excellent trade-off between attack detection and false-positive rates. We further propose an improved scheme based on software defined networks that automatically prevents threats by analyzing only the first few packets of a flow. The proposal promptly and efficiently blocks threats, is robust, and can scale up, even when the attacker employs spoofed IP.

KEYWORDS

big data, machine learning, stream processing, threat detection

1 | INTRODUCTION

Communication networks transfer an increasing amount data due to billions of users and devices, creating several challenges for threat detection.¹⁻³ Moreover, denial of service (DoS) attacks already reached more than 1.35 terabits per second⁴ and the longest in 2016 lasted 197 h.^{5,6} Network providers have to handle millions of data streams in real-time,⁷ and attacks can attempt to hide in this deluge of information; therefore, detecting network attacks is actually a “big data” problem and requires modern processing frameworks. Besides, the time it takes to detect an attack is essential, and if detection takes too long, irreparable damages will occur. Current security systems are not effective, since 85% of network intrusions are detected weeks after they had happened,⁸ with an average detection time of 206 days. To address the long time reaction to attacks, in this article, we leverage recently developed “big data” frameworks for real-time *stream processing*, and show their effectiveness for detecting network threats by providing *fast analysis of large and diverse* network data. We achieve high traffic processing rates, enabling the implementation of machine learning algorithms to detect both known and unknown attacks shortly after they occur.

In particular, we propose and implement a fast and accurate threat detection and prevention architecture that combines several big data open-source platforms for batch and stream processing. We propose a combination of (1) conventional *batch* processing over a historical database, with (2) real-time *stream* processing analysis to detect threats in real-time. We analyze the performance of three popular open-source stream processing platforms to choose the most effective one to detect threats as fast as possible.

For stream processing, we study two ways of combining the packets into flows, henceforward defined as a sequence of packets from the same source IP to the same destination IP within a time window. The first approach collects all packets in a time window, and the second approach uses only the first packets of each flow, that is, we periodically analyze the first few packets of each flow. Moreover, we propose a scheme based on software defined networks (SDNs) and the analysis of the first few packets of a flow to implement our threat prevention architecture, to block threats with spoofed IP addresses. Also, the off-path analysis machines only analyze the first packets of each flow per period, being immune to flooding attacks.

To evaluate our proposal, we create a dataset with labeled classes for the architecture evaluation, containing normal network usage and several attacks. Besides, we use another dataset with real-world broadband user data from one of the most significant network operators in Brazil. We implement five detection methods, three supervised classification methods to detect known threats, and two unsupervised anomaly detection methods to detect zero-day attacks and unknown threats. The results show a high accuracy for known threats, higher than 95%, and an excellent trade-off between false positives and threat detection in the anomaly methods.

In summary, the proposed threat detection and prevention architecture as a whole presents several advantages, such as (i) effective threat detection for both known and unknown threats; (ii) fast counter-measures against threats, since it relies on stream processing and the rapid analysis of the first packets of a flow, without having to wait to the end of the flow; (iii) robust against flooding attacks and with a high potential for scalability since the analysis machine only receives few packets per flow; and (iv) effective in the threat network traffic block, since, with the use of SDNs, the architecture blocks the malicious traffic as near of its source as possible, even on scenarios with spoofed source IP addresses.

We organize the rest of this article as follows. Section 2 discusses related work. Our proposed stream processing algorithms are presented in Section 3. Both datasets to evaluate the architecture are presented in Section 4. In Section 5, we discuss how to select the various parameters of our architecture and show the results of the different threat detection methods. Section 6 presents a resilient SDN-enabled architecture that can detect attacks accurately while surviving attempts by the attacker to obfuscate its activities and which also tries to attack our detection tools. Finally, Section 7 concludes the work.

2 | RELATED WORK

There are different approaches to detect DoS attacks. One of them is to accomplished statistical measures and trace back to the origins of the attack. The trace back procedure, however, requires to store information in the routers or to include information in every packet, which are complex operations.^{9,10} Our work is related to previous efforts on using *machine learning* and *big data* frameworks for detecting intrusions, and on SDNs to dynamically respond to detected threats.

Machine learning techniques can be either supervised or unsupervised, depending on whether the dataset is labeled or not. Algorithms that use supervised learning include neural networks, decision trees, and support vector machines (SVMs).¹¹ These types of algorithms have been increasingly applied to computer security applications; for example, Li et al. combine a pattern matching technique, dynamic time warping (DTW), and SVM to generate intrusion detection rules.¹² The method is evaluated with the traditional KDD dataset to classify between DoS, probe, and remote to local (R2L) attacks. Supervised learning, however, is trained with known attacks, and these algorithms sometimes have a hard time identifying previously unseen attacks. In unsupervised learning, the training data does not have labels, and the algorithms generally learn a representation of the patterns of the data. These algorithms are generally useful for anomaly detection; for example, Lakhina et al.¹³ propose the use of sample entropy for anomaly detection, and they show that this metric combined for source and destination IPs and ports, together with volume analysis can detect multiple sources of anomalies. Common methods to detect outliers apply principal component analysis (PCA) and independent component analysis (ICA) when data are assumed to follow a non-Gaussian distribution. Fernandes et al.¹⁴ use PCA combined with ant colony optimization for clustering, to perform profile-based anomaly detection. Palmieri et al.¹⁵ use ICA to separate the network traffic from their different sources. Then, decision trees are applied to perform binary classification. However, PCA and ICA are both sensitive to noise when used in anomaly detection.¹⁶ Amaral et al.¹⁷ proposed an improved approach using Tsallis entropy to detect anomalous traffic. The authors this technique to real and simulated traffic, but it only considers six features for anomaly detection. The authors use a graph representation of network features, allowing deep inspection of IP flows. In contrast with our parameter adaptation solution, this system needs parameters defined by the network administrator for its execution.

Chellammal and Malarchevi¹⁸ proposed an architecture that uses conventional algorithms in parallel and a feature selection mechanism to reduce data size. Nevertheless, the number of base learners used in the training phase and if they use the same or different base learns are application dependent, which makes complex the proposal deployment.

Villar-Rodriguez et al.¹⁹ propose the use of SVM to detect identity theft in social networks. The authors monitor user profiles based on connection time information. SVM classifies the legitimate user and the attackers' profiles. Li et al.²⁰ propose an approach to anomaly detection in traffic monitoring. The authors use PCA over the random forest machine-learning algorithm to identify the most important features. In our proposal, we improve the sample entropy metric by employing a time series that takes into account 26.

In addition to machine learning, the use of *big data* frameworks has also become popular in security applications. Bostani and Sheikan²¹ propose an unsupervised anomaly detection algorithm using MapReduce. The work uses the optimum-path forest (OPF) algorithm to project clustering models and detect anomalous behaviors. This article only focuses on two specific Internet of Things (IoT) attacks, sinkhole, and selective-forwarding,

disregarding unknown threats or zero-day attacks. Singh et al.²² proposed a similar platform in which the use of big data analytics leveraging Hadoop,* to detect peer-to-peer botnets. The offline analysis execution is the main limitation of this approach. BigFlow²³ is another proposal that employs stream processing for monitoring high-speed networks. The authors employ several classifiers to determine if a networking event is suspicious. Then, stream learning allows incremental model updates with new knowledge for the system. Even if the proposal can analyze events in a high-throughput, once a threat is detected, there is no protection system. Twister2 is a promising stream processing framework that presents a high performance and low latency, outperforming Apache Storm and Apache Flink.^{24,25} Nevertheless, at the moment, Twister2 only supports two machine learning algorithms.

In a previous work,²⁶ we present the tool CATRACA (scalable traffic classifier and analyzer), which implements an intrusion detection system using stream processing and machine learning. Nevertheless, in this article, we focus on providing real-time stream analytic to improve the time to detect attacks. Another related topic when using machine learning for security is how to evaluate the performance of these methods. Cardenas et al.²⁷ discuss the tradeoff between precision and false alarms and propose intrusion detection operating characteristic (IDOC) curves to evaluate intrusion detection problems. A classical dataset used for evaluating proposals is the KDD99 dataset,²⁸ a dataset created two decades ago.

A number of other works also create their dataset to evaluate their proposals.^{29,30} Sangkatsanee et al.²⁹ classify attacks using different algorithms, while Amini et al.³⁰ detect anomaly in real-time with unsupervised neural networks. These works, however, do not evaluate their proposals in scenarios with a large amount of traffic. On the contrary, we use two real traffic datasets (one with a large amount of traffic) with packet captures from our lab and one of the major network operators in Brazil.

The programmability of software defined networking is also receiving more attention in the security community. We proposed BroFlow,^{31,32} a combination of the Bro IDS and the OpenFlow POX controller. We implemented an algorithm to block a DoS attack by using Bro to send messages to the controller, to mitigate the attack automatically. Nevertheless, the implemented counter-measure blocks all traffic generated by the identified source IP. Therefore, the system is ineffective under IP address spoofing. In his Ph.D. thesis, the co-author Andreoni Lopez³³ proposes a threat detection system that presents the following contributions: a virtualized network function in an open source platform, a fast feature selection algorithm, a near-optimal placement of sensors, and a greedy algorithm that allocates on demand a sequence of virtual network functions. Lin et al.³⁴ extend the SDN architecture for traffic classification and intrusion prevention. The proposal detects DoS attacks and malicious HTTP requests. Nevertheless, the work mainly uses SDN functionalities for load balancing, and it does address automatic attack detection. OPcloudSec is another proposal that integrates SDN and cloud computing for security management.³⁵ The authors use a deep learning approach based on Deep Belief Networks for DDoS detection. The authors do not discuss the requirement of parameter tuning and disregard the effect of IP Spoofing when detecting DDoS. IP spoofing could easily overload routing tables in SDNs. Vincentini et al.³⁶ use a similar approach to our proposal combining Apache Storm for streaming processing with Floodlight controller for SDN. In contrast to our work focused on network security, the authors use this combination for network management in a multi-tenant environment.

In a companion conference paper, we present preliminary results of an adaptive threat detection architecture that uses a honeypot and trains detection models in real-time.³⁷ Our proposal extends our previous work in several aspects: (i) we leverage big data stream processing frameworks, (ii) we propose fast machine-learning algorithms for the first few packets of each flow, and (iii) we use SDNs to respond to an attack promptly. Unlike the previously cited papers, we propose a specific architecture that uses the lambda architecture, allowing real-time stream processing analysis based on the support of a historical database. Our architecture prevents threats in a fast and scalable way by providing real-time, accurate detection of known and zero-day attacks through automated classification and anomaly detection methods.

3 | PROPOSED THREAT DETECTION ARCHITECTURE

Conventional *big data analytics* tools usually employ batch processing; however, batch processing produces high latency, with responses in the order of tens of seconds, which is unacceptable for critical applications requiring real-time processing, with fast responses within a second.³⁸ Unlike batch processing, stream processing techniques can analyze continuously generated data and provide real-time results. Their accuracy, however, might not be as good as batch methods that use more data. Both of these paradigms, batch, and stream, can be combined to exploit the benefits of each technique in the lambda architecture, which analyzes big data in a real-time and accurate manner.³⁹

Our proposed architecture, shown in Figure 1, is based on the lambda architecture because it combines both *batch* and *stream* processing methods. The lambda architecture has three layers: the stream processing layer, the batch-processing layer, and the service layer. The stream processing layer deals with the incoming data in real-time. The batch-processing layer analyzes a huge amount of stored data in a distributed way through techniques such as map-reduce. Finally, the service layer combines the obtained information of the two previous layers to provide an output composed of analytic data to the user. Therefore, the lambda architecture gives us the ability to analyze streaming data to obtain real-time results while at the same time complementing this analysis with historical data.

*Hadoop is the Apache foundation framework to store and process data using MapReduce processing model.

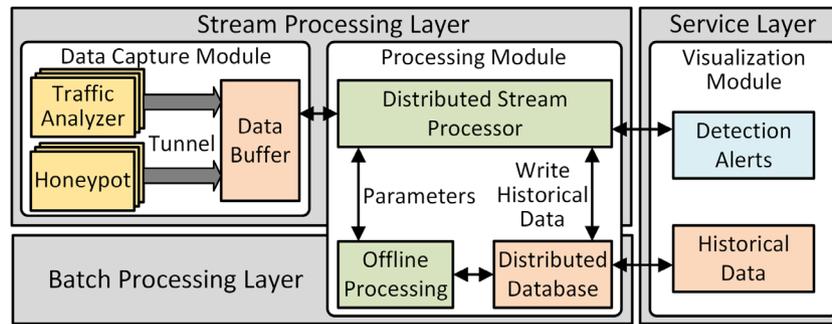


FIGURE 1 Proposed architecture for real-time stream processing and historical data batch processing composed of the following modules: (i) the data capture module gathers data, (ii) the processing module analyzes and stores it, and (iii) the visualization module displays analytic information results

We implement the traffic analyzer using `Bro`,[†] a real-time network security monitor with its network-oriented programming language that allows us to create programs to make flow abstraction easier. In our proposed architecture, `Bro` identifies network flows, extracts information from the packets, and groups them in time windows.

Since network usage varies over time, captured data arrives at different rates. Hence, an overload may occur while processing the streams, resulting in data loss. Apache `Kafka` receives the captured data, acting as a buffer to prevent overload at the analysis location. `Kafka` is a message broker that works as a publish/subscribe service and acts basically as a buffer to the processing tool, adapting different generation and processing rates. `Kafka` abstracts the message flow into topics. Producers then write their data into topics from which the consumers can access the information.

To analyze streaming data, we chose `Storm` over two other streaming platforms based on a performance evaluation we develop in Section 3.1. `Storm` offers a distributed fault-tolerant stream processing framework. In addition, `Storm` processes data in memory, ensuring low latency in real-time. `Storm` processes the streams in a direct acyclical graph (DAG) composed of input elements, called *Spouts*, and processing elements, called *Bolts*. The application can define the parallelism of the *Spouts* and *Bolts* in a way that multiple stream samples can be processed simultaneously.

Once we obtain analytics from the streams, we store the results in a dynamic database for historical analysis. We use a distributed database to achieve better resilience. Our architecture uses Apache's `HBase`, a fault-tolerant database that can store large quantities of sparse data. We calculated offline the parameters with historical data to improve the accuracy and adaptiveness of our threat detection. Then, we calibrate the processing model for real-time threat detection. The architecture has an adaptive characteristic, because the parameters are periodically updated, adapting to new network use patterns.

3.1 | Stream processing performance evaluation

We evaluate three open-source streaming platforms with our threat detection application to choose the one that can process data at higher rates because one of our main concerns is to detect threats as fast as possible. We compare `Storm` against `Spark Streaming` and `Flink`.⁴⁰ We use our threat detection application as a benchmark to measure stream-processor performance. The attack-detection methods in our application are further presented in Section 5. The experiment evaluates the performance of the platforms in terms of processing throughput. The dataset is injected into the platforms in its totality and replicated as many times as necessary. We measure the consumption of messages and the processing rate of each platform. We also vary the parallelism parameter, which represents the total number of cores available for the cluster to process samples in parallel. Figure 2(A) shows the performance results of this experiment. `Storm` shows a higher throughput when compared to the alternatives. Figure 2(A) also shows that `Storm` process up to 15 million samples per minute with our threat detection application, which gives about 4 μ s of detection time, allowing defense strategies and significantly decreasing the risks.

In addition to this experiment, we also use another benchmark that counts the number of times each word appears in a text, using a dataset that contains more than 5,000,000 tweets.⁴¹ All three platforms offer the word-count application as an example. Therefore, we show this result to get an unbiased comparison that is not affected by our implementation. Figure 2(B) shows the performance behavior of the three systems under a word-count program. Once again, `Storm` has a better performance and, therefore, is the adequate platform for our threat detection architecture.

[†]The `Bro` project was renamed to `Zeek` project: <https://www.zeek.org/>

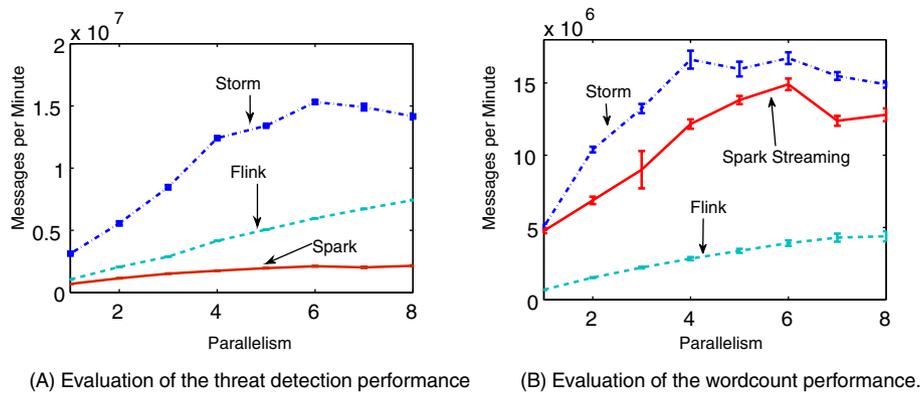


FIGURE 2 Throughput results of the platforms in terms of number of messages processed per minute in function of the task parallelism

4 | SECURITY DATASETS

There are few openly available network security datasets to evaluate defense mechanisms due to the privacy and security concerns with sharing real-world data. DARPA⁴² is the first effort for creating open network intrusion data and its follow-up KDD 99²⁸ dataset. The DARPA dataset includes TCP/IP traffic and operating system data collected from a simulated computer network. While collecting the data, attacks were also simulated and labeled in the dataset. The KDD 99 dataset consists of a selection and grouping of DARPA features to facilitate the application of machine learning algorithms. These datasets, however, consider a network simulation that is a limitation⁴³ because it introduces artifices in the training and testing dataset. Moreover, these datasets are two decades old and do not represent current threats.⁴⁴

In this work, we evaluate our threat detection proposal using two different datasets, both containing real-world traffic. One dataset contains traffic from a Brazilian Internet Service Provider, and the other contains real traffic from our lab. The use of two different datasets shows that the proposed architecture and its detection methods work well, even considering distinct scenarios. Moreover, we propose two modes of combining the packets into flows, both of them considering a flow as a sequence of packets from the same IP source to the same IP destination during a period. In the former, we gather all packets in a fixed-length time window. We determine the length of this window further in this article, based on the accuracy of the classification algorithms. Each flow has 26 features, generated by the TCP/IP header data. The main features are TCP, UDP, and ICMP packet rates; source and destination port number; the number of each TCP flag; average and variance of inter-packet arriving time; average and variance of flow packet length; among others. The second mode to define flows consists of extracting the features from the first few packets of each flow. The intuition behind this approach is the well-defined initial behavior of most applications, which leads to proper classification. Once again, we determine the number of packets to be considered based on the accuracy of the classification methods.

4.1 | Network operator dataset

Real-world information from 373 residential broadband users from the city of Rio de Janeiro for one week composes the Network Operator (NetOp) dataset. We anonymize the network traffic for privacy considerations. An intrusion detection system (IDS) filtered the traffic. We analyzed the logs from this IDS, and the proportion of traffic filtered out was around 15%. Since we obtained the data filtered by an IDS, we added real botnet malicious traffic captured in the work of Garcia et al.⁴⁵ to detect these threats and evaluate our threat detection architecture. The botnet data have 13 different scenarios of malware infection. These attacks are real and were not performed by the authors since they infected the machines with real malware. In the combined dataset, we keep 15% threat traffic proportion.

4.2 | GTA/UFRJ dataset

Another contribution of this work is the creation of a dataset with real network traffic to evaluate the proposal.[‡] The dataset has around 95 GB of packet capture raw data in computers from our lab, GTA at Federal University of Rio de Janeiro. We added to the normal traffic, real network threats, including seven types of DoS and nine types of network probes. The analysis of packet header information detects two threat classes: DoS attacks

[‡]The dataset can be obtained by emailing the authors.

and probe. Therefore, we elaborate on the dataset with several attacks from both these classes. The DoS attacks are *ICMP flood*, *land*, *nestea*, *smurf*, *SYN flood*, *teardrop*, and *UDP flood*. The different types of probes in the dataset are *TCP SYN scan*, *TCP connect scan*, *SCTP INIT scan*, *Null scan*, *FIN scan*, *Xmas scan*, *TCP ACK scan*, *TCP Window scan*, and *TCP Maimon scan*. We launch the attacks using tools from *Kali Linux*. These attacks were labeled in the dataset by origin and destination IP filters, separating the traffic belonging to the attack machines from the normal lab network usage.

5 | AUTOMATIC THREAT DETECTION

Our proposed real-time architecture relies on machine learning algorithms to perform automatic threat classification. In Sections 5.1–5.3, we present the classification algorithms implemented to evaluate the architecture. We selected these algorithms because they are among the most popular for network security.^{11,46} In all methods, we perform the training with 70% of the dataset and the test with the remaining 30%. During the training phase, we perform tenfold cross-validation to avoid overfitting the data.

5.1 | Decision trees

In the decision tree algorithm, leaves represent the final class, and branches represent conditions based on the value of one of the input features. During the training part, the C4.5 algorithm determines a tree-like classification structure, based on the information entropy of each feature. The real-time implementation of the decision tree consists of if-then-else rules that generate the tree-like structure previously calculated. The results are presented in Section 5.6, along with the other algorithms results. During the training phase, we set the following parameters: maximum number of splits as the number of samples minus one; minimum leaf size as one; and minimum number of branch node observations as ten observations. We chose these parameters because they ensured an excellent trade-off between the number of leaves and the tree depth. A deep tree with many leaves may result in high training accuracy but does not perform well in independent test sets.

5.2 | Artificial neural networks

Artificial neural networks were originally inspired by the human brain, in which each neuron performs a small part of the overall processing, transferring the output to the next neuron, and achieving results from the combination of these subtasks. In classification neural networks, the final output represents a degree of membership for each class, and the output with the highest membership degree determines the predicted class. The training phase tunes the neural network adjusting the weight vectors Θ . These vectors determine the weight of each neuron connection. During the training phase, the input vectors are mapped into a predicted output vector and compared to the real output. The prediction errors are then minimized by the back-propagation algorithm, taking into account the error value induced by each parameter.

In the classification phase, each neural network layer computes the following equations:

$$z_{(i+1)} = \Theta_{(i)} a_{(i)} \quad (1)$$

$$a_{(i+1)} = g(z_{(i+1)}) \quad (2)$$

$$g(z) = \frac{1}{1 + e^{-z}} \quad (3)$$

where $a_{(i)}$ is the vector that determines the output of layer i , $\Theta_{(i)}$ is the weight vector from layer i to layer $i + 1$, and $a_{(i+1)}$ is the output of layer $i + 1$. The function $g(z)$ is the *Sigmoid* function that plays an important role in the classification. For high values of z , $g(z)$ returns one and for low values $g(z)$ returns zero. Therefore, the output layer gives the degree of membership of each class, between zero and one. The trained neural network consisted of three layers, the input layer, the hidden layer, and the output layer. The input layer has 26 neurons and the output two since we classify each flow as a threat or normal. The hidden-layer size is equal to 10 neurons. We chose the regularization parameter as one, to prevent data overfit.

5.3 | Support vector machines

SVMs are binary classifiers based on the concept of a hyper-plane in a multidimensional space that splits different classes. An iterative algorithm minimizes an error function, finding the best hyper-plane separation. A kernel function defines this hyper-plane. This way, SVM finds the hyper-plane of maximum margin, that is, the hyper-plane with the biggest distance possible to both classes.

The classifier score for a sample x is the distance from x to the decision boundaries, which go from $-\infty$ to $+\infty$. The classifier score is given by:

$$f(x) = \sum_{j=1}^n \alpha_j y_j G(x_j, x) + b, \quad (4)$$

where $(\alpha_1, \dots, \alpha_n)$ are the estimated parameters of SVM, and $G(x_j, x)$ is the used kernel. In this work, the kernel is linear, that is, $G(x_j, x) = x_j^T x$, which provides good performance with the minimum number of input parameters. We set the cost parameter as one to control the margin-violation penalty on observations, and to reduce overfitting. The initial coefficients, α_i , were randomly initialized.

5.4 | Flow parameter tuning

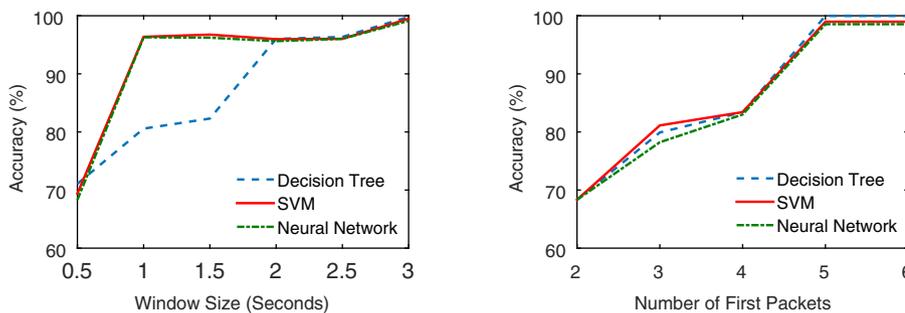
We have to determine two critical parameters for our threat detection architecture: the window-time size for flows, and the number of initial packets needed to characterize a flow accurately. Small window-time sizes provide faster threat detection but can compromise accuracy. Similarly, the fewer packets required to classify a flow, the quicker our algorithms can reach a decision. Figure 3 shows low-accuracy results for 0.5 second window size and two packets. This accuracy result, however, improves as the architecture gathers more information about the flows. In our previous work,³⁷ we suggested that a network flow sampling as a sliding window of 2 s duration is the best trade-off between classification accuracy and decision latency.

We use the GTA/UFRJ dataset and the classification algorithms implemented to choose these parameters. Figure 3(A) shows the accuracy result for the approach with all packets gathered in different time windows. Both neural networks and SVMs provided good detection results, starting with a window size of just 1 s. At the same time, decision trees only achieved similar results with window sizes of at least 2 s. We choose a window size of just 1 s because that is enough information to classify the samples with the shortest possible time correctly.

We now turn our attention to selecting the number of packets to analyze. Figure 3(B) shows the accuracy for all three algorithms using the approach that analyzes the first few packets. Usually, the behavior of applications and threats is well defined early on, and as the result shows, five packets are enough to obtain high accuracy, near 99%. The first five packets approach shows higher accuracy than the method that analyzes all packets for one-second window size. This approach presents other advantages, such as shorter time to extract flow characteristics, greater robustness, and better scalability. Furthermore, it is robust to flooding attacks since there is no need to process a large number of packets for each flow.

5.5 | Feature selection and PCA

We perform dimensionality reduction to improve the efficiency of the proposed architecture for real-time threat detection. The aim is to improve the throughput, eliminating irrelevant features from the threat detection procedure. We also study possible correlations between two or more features, which can be combined into only one feature, reducing processing time. We perform dimensionality reduction with PCA, which transforms a group of possibly correlated variables into a group of weakly correlated variables that lie in orthogonal planes. This transformation takes into account the eigenvalues in a way that the component associated with the biggest eigenvalue represents a more significant data variance. The algorithm sorts in variance order the other components to obtain the resulting matrix. We keep the components related to the higher eigenvalues because they have more relevant information, and we remove the components associated with low eigenvalues, reducing data dimensionality, and improving the processing time. It is important to remark that PCA does not consider the class label in the dataset, and therefore, can be used in both supervised and unsupervised learning.



(A) Accuracy for different flow time window size. (B) Accuracy for different number of first packets per flow.

FIGURE 3 Accuracy of decision tree, SVM, and neural network algorithms for the two flow approaches, all packets in a fixed time window size and first few packets of a flow

For both ways of combining packets into flows, the sum of the higher eight eigenvalues represents more than 95% of the total amount, as shown in Figure 4. In other words, the first eight features from the data calculated by the PCA linear transformation represent 95% of the total variance. Therefore, we select these eight components and discard the others that represent less than 5% of the total data variance, improving the processing time, which is critical in real-time applications.

5.6 | Threat classification results

Tables 1 and 2 show the accuracy comparison for the first five packets of each flow, and the one-second time window approaches for both the GTA/UFRJ and NetOp datasets. We trained the classification algorithm using all attacks in the dataset with the same label, therefore, our problem is the binary classification between threats and legitimate traffic. The results show that using the first five packets in the time window provides higher accuracy than the one-second time window for SVM and neural network classifiers. The accuracy improvement is due to the behavior of applications and threats, which are defined in the beginning, usually when the negotiation phase of applications happens. Both SVMs and Neural Networks perform well in all scenarios. As shown in Tables 1 and 2, the accuracy for both these algorithms is higher than 95% for all scenarios, ensuring the efficiency of the proposed architecture to detect known threats.

Using the first five packets of each flow gives us better results than using all the packets in a one-second time window. Our proposal randomizes the time for checking the packets, and therefore, the attackers cannot easily simulate a legitimate use for a fixed time interval and then engage an attack.

We further detail precision and recall for both legitimate and malicious traffic in Table 3. Precision and recall are standard metrics for evaluating the results for a specific class. They are recommended metrics for evaluating intrusion detection systems (they are similar to IDOC curves²⁷), as they provide more meaningful results than ROC curves and the false positive versus false-negative rates.²⁷ Since the primary goal of our threat detection architecture is to detect threats and trigger counter-measures, we show these metrics for the threat and normal classes, to indicate the number of true and false alerts that would trigger counter-measures.

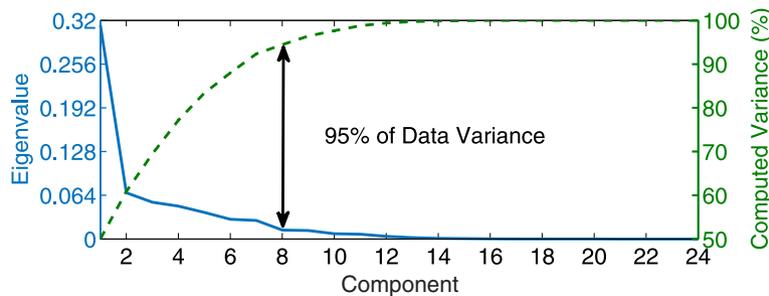


FIGURE 4 Eigenvalue for each flow feature. The eigenvalue associated to each of the transformed features is proportional to the data variance. The eight highest principal components represent 95% of the total data variance

TABLE 1 Accuracy comparison for the three classification methods in the GTA dataset

	First packets of flow	One second window
Decision tree	99.9%	80.6%
Neural network	99.0%	96.0%
SVM	98.6%	96.3%

TABLE 2 Accuracy comparison for the three classification methods in the NetOp dataset

	First packets of flow	One second window
Decision tree	86.3%	92.8%
Neural network	95.3%	95.1%
SVM	96.1%	95.8%

TABLE 3 Threat classification summary for all algorithms and datasets

Algorithm	Dataset	Threat as positive	
		Precision	Recall
Decision tree	GTA 1s	80.1%	94.6%
	GTA 5p	99.9%	99.9%
	NetOp 1s	97.3%	94.8%
	NetOp 5p	99.1%	85.6%
Neural network	GTA 1s	94.8%	98.8%
	GTA 5p	98.6%	99.8%
	NetOp 1s	97.5%	97.2%
	NetOp 5p	98.2%	96.6%
SVM	GTA 1s	95.4%	98.8%
	GTA 5p	98.0%	99.7%
	NetOp 1s	96.5%	98.9%
	NetOp 5p	99.6%	96.1%

Note: We show precision and recall to evaluate the methods. The 1s stands for the one second time window approach for combining flows and the 5p stands for the first five packets in the period.

For any given class, the precision is the fraction of correctly classified samples over all samples predicted to belong to such class. At the same time, the recall is the fraction of correctly classified samples over all samples that belong to that class. As Table 3 shows, we achieve good results, usually above 90%, on both precision and recall for the normal class. Consequently, this ensures a low level of false alerts that would result in legitimate traffic blocking. Regarding the threat class, for the NetOp dataset, the precision results are lower, because the dataset has 85% of legitimate samples. Therefore, when evaluating the absolute number of normal samples classified as threats, they have a negative impact on the precision. Nevertheless, as shown in the results for the normal class, they do not result in many false alerts. Taking this into consideration, the most interesting measure to evaluate for the threat class is the recall that measures the threat-detection rate, that is, the number of correctly classified threats among all real threats in the datasets. The results show that the first five packets of each flow present a higher threat recall, therefore, a higher threat detection rate. Moreover, SVM is the most stable method, since, for all scenarios, the threat recall is above 87%.

5.7 | Anomaly detection by normal distribution

So far, we have used supervised-learning algorithms to detect previously seen attacks; however, protection against unknown attacks is essential to have a higher level of security in computer networks. Zero-day attacks are hard to detect since there is no previous information about the attack. Anomaly detection is capable of discovering these new attacks by identifying abnormal flows.

In the next two subsections, we present and compare two anomaly detection algorithms; (1) anomaly detection by distance to a normal distribution, and (2) an anomaly detection algorithm based on entropy metrics.

In the first method, we detect anomalies through the mean and variance from each feature of the normal samples of the training dataset. We identify anomalies when the distance from the sample feature to the mean is greater than a threshold times the variance in at least one of the features. We analyze the eight PCA transformed features.

The real-time implementation requires anomaly detection as the streaming data is arriving. The anomaly is detected if at least one of the following conditions is true for at least one feature j , taking into account the means μ_j and the variances σ_j^2 calculated in training:

$$X_j > \mu_j + \text{threshold} * \sigma_j^2 \quad (5)$$

$$X_j < \mu_j - \text{threshold} * \sigma_j^2 \quad (6)$$

The proposed architecture allows real-time anomaly training. Consequently, the algorithm becomes adaptive, which is fundamental for anomaly detection, since the network behavior may change in time. Therefore, when a new sample arrives and it is not detected as an anomaly by conditions (5) and (6), the parameters μ_j and σ_j^2 of each feature are updated, considering this new sample. The parameters of a normal distribution are expressed by:

$$\mu_j = \frac{1}{N} \sum_{i=1}^N X_j \quad (7)$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (X_j - \mu_j)^2 \quad (8)$$

Therefore, current values of the sum and the total of samples N are stored and incremented when a new sample arrives, considering each feature X_j . As a consequence, samples considered legitimate to update the normal distribution parameters, ensuring adaptability.

Figure 5 shows the normal distribution anomaly detection and false-positive rates for the NetOp dataset. For the first packets of each flow approach, with a threshold value of 2.1, we achieve a remarkably low false positive rate of 0.7%, maintaining a good threat detection rate of 87.7%. For the one-second window approach, considering a threshold of 3, we obtain 6.7% and 88% false-positive and threat detection rates, respectively.

Figure 6 shows the results for the GTA/UFRJ dataset, considering different threshold values. We obtain the false positive rate using the normal test dataset and the threat detection rate using all the threats in the dataset. When choosing a low threshold value, the proposal detects almost all threats. Nonetheless, this comes with the cost of a high false-positive rate. On the other hand, a high threshold value results in fewer false-positives, but also a lower threat detection rate. With a threshold value of two, the false positive rate is 5.6%, and the detection rate 96.4% for the one second time window approach. With a threshold of 2.7, these rates are 7.9% and 88.6% for the first five packets of each flow approach. The value for the threshold parameter depends on the application and considers a trade-off between the false-positive and attack detection rates.

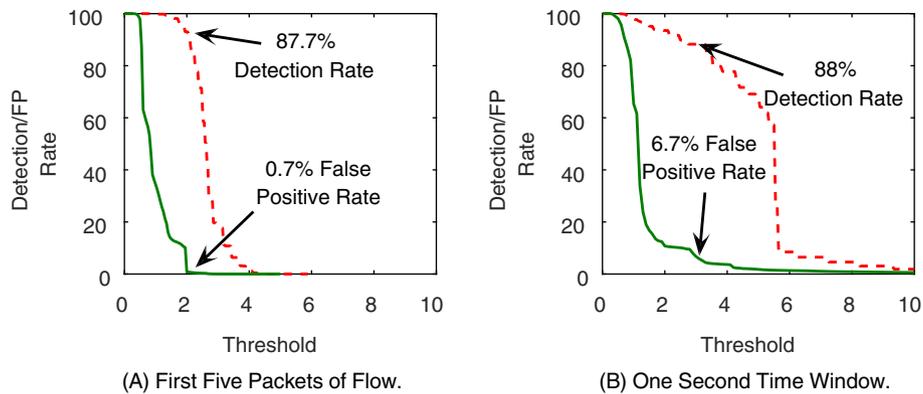


FIGURE 5 False positive and attack detection rates for the NetOp dataset according to the threshold

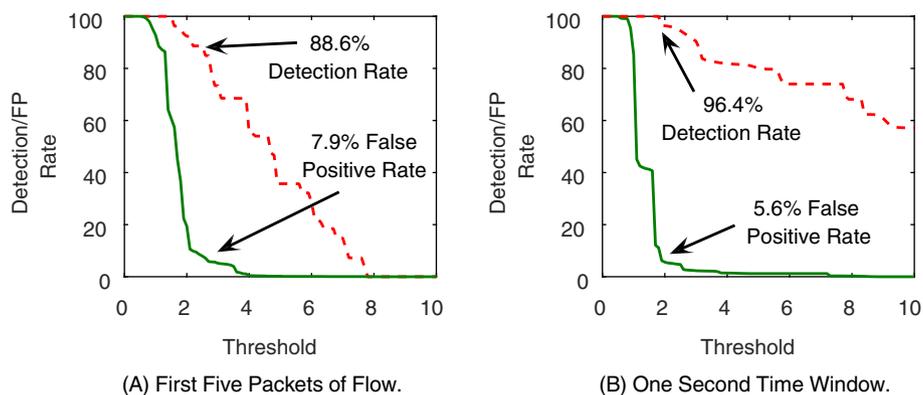


FIGURE 6 False-positive and attack detection rates for the GTA/UFRJ dataset according to the threshold. The lower the threshold, more attacks are detected, but also higher is the false positive rate

5.8 | Anomaly detection by entropy time series

We also implement an anomaly detection algorithm by analyzing the entropy value of a sliding window of flows. The sample entropy indicates the degree of concentration or dispersion of a characteristic. It is calculated as follows:

$$H(X) = -\sum_{i=1}^N \left(\frac{n_i}{S}\right) \log_2 \left(\frac{n_i}{S}\right) \tag{9}$$

where S is the total number of observations and n_i is the number of observations within the range i of values, and N is the number of ranges. When one range concentrate all values, $H(X)$ is equal to zero and when each value is in a different range i , the value of $H(X)$ is $\log_2(N)$. Therefore, given a series of observations X , the sample entropy summarizes the level of concentration in one single value.

To detect anomalies using the sample entropy value, we define a sliding window of 40 flows and calculate the value of $H(X)$ for each of these windows, generating the time series. In the training phase, we calculate the histogram of all normal samples of the training dataset and determine 30 ranges of equally distributed entropy values. Another parameter specified in the training phase is the range of most samples. We observed that the normal traffic entropy values tend to concentrate together and that usually, the most frequent value is in the middle of this concentration. Thus, we propose to detect anomalies by defining a threshold that limits the accepted distance of the entropy $H(X)$ to the most frequent range. The architecture can update the most frequent value online, adapting to different network behaviors.

Figure 7 shows the results for different threshold values considering the GTA/UFRJ dataset. For the first five packets of each flow approach, with a threshold of 0.8, the threat detection rate is 92.1%, and the false-positive rate is 7%. For the one-second window and the threshold value of 1.3, these values are 86.8% and 2.8%. Once again, we can determine this threshold, considering the trade-off between threat detection and false-positive rates.

Figure 8 shows the entropy anomaly detection results for the NetOp dataset. For the first packets of each flow approach, the threat detection rate is very high, 91.8%, and the false-positive rate is low, 1.5%. The result for the one second time window is as good as the other approach, with

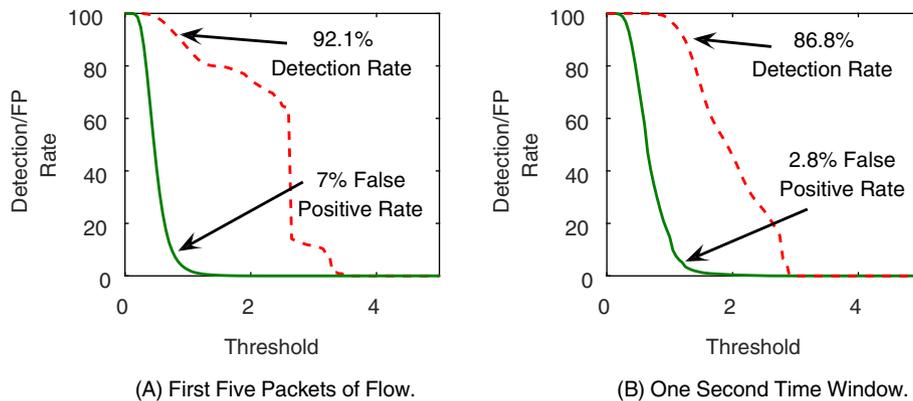


FIGURE 7 False positive and attack detection rates for the GTA/UFRJ dataset according to the entropy threshold. The threshold represents the distance to the range that has the most entropy samples

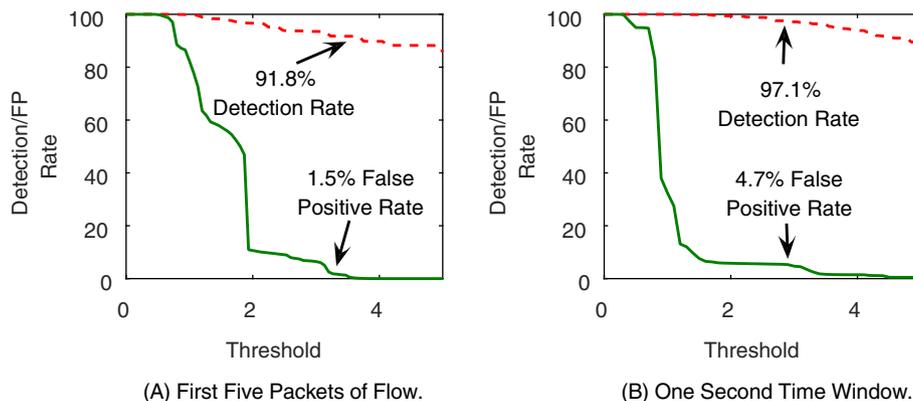


FIGURE 8 False positive and attack detection rates for the NetOp dataset according to the entropy threshold. The threshold represents the distance to the range that has the most entropy samples

false-positive and detection rates of, respectively, 4.7% and 97.1% with a threshold of 3. The results for this method show an excellent trade-off in all scenarios.

6 | THE THREAT PREVENTION ARCHITECTURE

We are now ready to have all our algorithms in a combined threat prevention architecture that analyzes traffic online, in batch mode, and that, also, provides fast attack response algorithms using the programmability of SDNs. We classify intrusion detection systems in two modes: on-path and off-path. In the on-path mode, we accomplish the analysis in the path from the source to the destination. In the off-path mode, the packets are mirrored and sent to an investigation that will determine whether the packet is part of a threat. The great advantage of the on-path mode is the ability to block traffic actively. This approach, however, introduces latency problems, considering that we add the processing time to analyze before the forwarding procedure. Moreover, the off-path mode allows gathering information from other sources. Thus, the system can correlate essential data to improve detection accuracy. On the other hand, the off-path mode needs an additional infrastructure to carry both traffic mirroring, and threat blocking.

We perform threat detection based on information from the first five packets of each flow (as discussed in the previous section). Whenever an undefined flow arrives at an SDN switch, it forwards a message to the controller, which then installs a rule with two actions, the first one forwarding the flow to its destination and the second one replicating the flow to an analysis machine. Meanwhile, the analysis machine keeps track of the number of packets for each flow that it is analyzing. When that number reaches five packets, this machine sends the flow information to the threat detection application and also a message warning to the controller, informing the accomplishment of packets analysis of a particular flow. We define flow as all packets from the same IP source and to IP destination, and the analysis machine extracts 26 flow features and publishes it in the `Kafka` message broker so that the processing core of the detection architecture gets the information and determines whether it is a threat.

Once the controller receives the message indicating the flow analysis completion, it lists all flows with the source and destination IP of the specified flow and removes the mirroring rule. The controller keeps the routing action to the flow destination, which preserves network operation. Thus, after analysis, the flow is no longer forwarded to the analysis machine, making the process robust because the analysis machine is protected against flooding attacks and has a greater ability to analyze flows when compared to the alternative of analyzing all packets. It is important to remark that the controller installs and removes the mirroring action periodically to increase network security. Therefore, after removing the mirroring action, it reinstalls again after the flow timeout set in the controller, which forces the periodical check of the flows. This timeout is randomly chosen within a range of values, to avoid an attacker that tries to bypass our system by sending legitimate flows during the beginning of the analysis.

Figure 9 shows an example network analyzed and protected by the proposed threat detection and prevention architecture. All SDN switches have an interface to which traffic is duplicated and sent to an analysis machine. The controller is responsible for installing the mirroring rule and for removing it when the machine accomplished the reception of the five required packets. The capture module, composed of the analysis machines, then forwards the flow characteristics to the detection application, which utilizes machine learning algorithms in conjunction with stream processing in real-time to do its analysis. If a threat is detected, the detection application sends an alert to the controller, which then blocks the attacker's source IP in all switches.

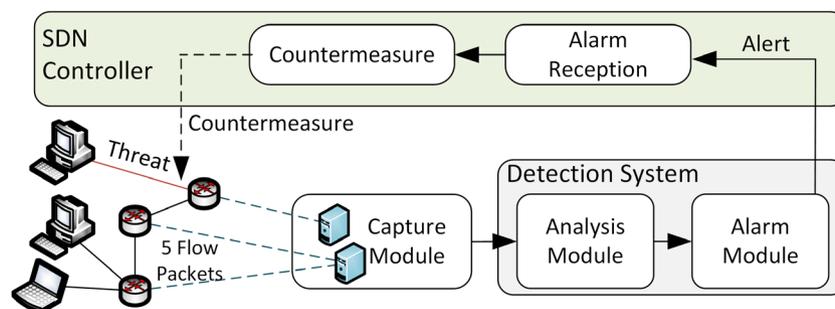


FIGURE 9 Example of network monitored by the proposed architecture. The network controller installs the mirroring rules for the capture module that sends traffic information to the detection system. Through alerts from the detection system, the network controller can block the threats

6.1 | Defense strategy against attacks with spoofed source IP

The source IP blocking rule is ineffective against attacks that falsify the source IP address of their packets to deceive defense systems. Attacks with spoofed source IP are even more critical in SDNs. Besides the attack damage, they overload the controller to set new flows every time the source IP changes, also generating a large number of entries in the switches flow table.

We propose a strategy against spoofed IP attacks based on a sequence of alerts and marking the path of flows. The intuition behind this concept is as follows: if the controller installed a blocking rule against an attack, and even then, an alarm arrived soon after, it may indicate that the blocking rule was not effective. Therefore, the controller keeps track of the time between the reception of the alert and notices when two alerts arrive in a short time. When this happens, the controller begins to suspect that the packet belonging to the attack has a spoofed IP. Due to this suspicion, the controller starts to map the path of all network flows for a specified period. This mapping is possible since, in SDNs, the controller has a global network view and knows the entire network topology. Then, if a third warning reaches the controller, in addition to the traditional source IP blocking rule, the controller also finds out in which switch port the attacker traffic enters the network and blocks this port. Thus, the controller instantiates a blocking rule of the traffic from the attacker TCP port, in addition to the IP blocking rule. Here we block the traffic from the port, considering that the intruder is in the SDN network. Nevertheless, other actions are applicable, such as limiting the bandwidth of the suspected attacker or transferring the traffic to a honeypot. Therefore, if an attacker spoofs its source IP, after three warnings, the attack traffic is blocked. Asking the controller to monitor the path of all flows only occurs when two alerts happen in a short period. On the other hand, after a period without receiving any new alert, the process is undone, avoiding the waste of resources. Another consideration is that instead of an attack with spoofed IP, it could be a distributed attack. Against distributed attacks, it is essential to block all sources, and thus the blocking rules are effective when blocking traffic from both the inbound network interface and the source IP.

Algorithm 1. Defense strategy against attacks with spoofed source IP

input : Incoming Alerts from the Detection Methods

output: Blocked source IP or incoming port

Initialize threshold, status;

while *True* **do**

wait(alert);

block(alert.sourceIP);

switch *status* **do**

case *normal* **do**

if *alert.time – lastAlert.time < threshold* **then**

mapSourcePort(start);

status = monitoring;

end

end

case *monitoring* **do**

if *alert.time – lastAlert.time > threshold* **then**

mapSourcePort(stop);

status = normal;

else

block(sourcePort(alert.sourceIP));

end

end

end

lastAlert = alert;

end

Algorithm 1 summarizes the defense strategy against spoofed IP addresses. The controller has two states: normal and monitoring. In the normal state, there is no suspicion that an attack with spoofed IP is occurring. In the monitoring state, there are suspicious activities, and the controller maps the source port of every incoming flow, using its global network view. The controller changes the status based on the time difference between two consecutive alerts. When this difference is above a threshold, the controller either keeps the monitoring status or goes back to the

normal status. Similar behavior occurs when this difference is below the threshold, however, the controller changes to or keeps the monitoring status. An essential aspect of the proposed strategy is that we only map the flows under spoofed IP suspicion, and we return to the normal status to avoid a waste of resources. The first step of the algorithm once an alert arrives is to block its source, therefore, preventing the network against distributed threats by blocking each one of them. Here, we implemented the action of blocking the source port when detected a threat that probably is spoofing its IP. Other counter-measures, however, are applicable, such as reducing the bandwidth or forwarding to a honeypot or network filter.

6.2 | Traffic monitoring and threat blocking

In this article, we use the Future Internet Testbed with Security (FITS) platform that is a combination of Xen hypervisor for virtualization and OpenFlow for traffic forwarding.⁴⁷ Figure 10 shows the constructed topology for the experiments, which consists of three client machines and a server machine. The packet forwarding is accomplished by OpenvSwitch switches that are controlled by an application programmed in the POX controller. Furthermore, an analysis machine characterizes flows using `BR0`. The connection between switches and the analysis machine requires a generic routing encapsulation (GRE) tunnel, which encapsulates the packet and sets the address of the analysis machine as a destination. The analysis machine parses the packets and then analyzes them. We perform the experiments on an Intel Xeon X5690 server with 24 processing cores, each with a frequency of 3.47 GHz clock and with 48 GB of RAM.

The results of the first experiment, shown in Figure 11, aim to display the operation of the traffic mirroring rule and the subsequent end of this rule after the analysis of the first five packets. In this experiment, the three clients send traffic at a constant rate to the server machine. Figure 11(A) shows the traffic received by the server machine. The results show that the proposed scheme does not affect communication. Furthermore, Figure 11(B) shows the packet rate received by the analysis machine, which sends a message to the network controller after capturing the five packets needed to characterize the flow. Even though the sending rate is much higher, the analysis machine receives a few packets. The reason the analysis machine receives a little more than five packets is the time required to inform the controller to undo the mirroring rule. Besides, this figure

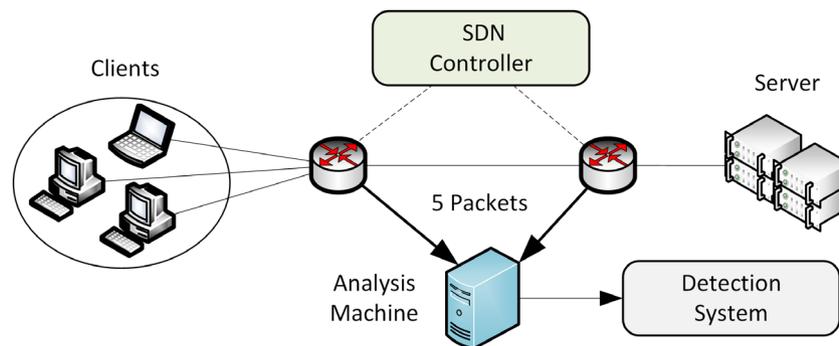
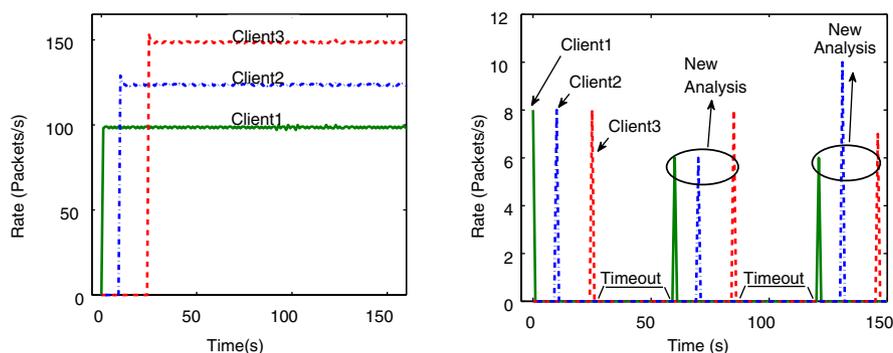


FIGURE 10 Experimental network topology used for the five packets scheme and the malicious traffic block



(A) Packet rate received by the server machine. (B) Packet rate received by the analysis machine.

FIGURE 11 Network operation when three machines communicate with the server. The server receives traffic without being affected by the proposed scheme, while the analysis machine only receives the first packets of each flow

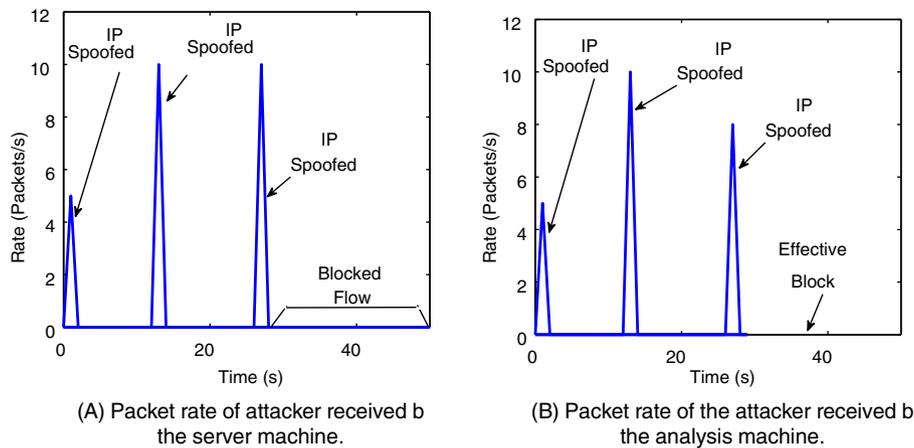


FIGURE 12 Architecture operation under an attack that spoofs the source IP address. Blocking rules cease to be effective when the IP is spoofed, however after the identification of the interface through which the attack enters the network, the attack is effectively blocked

shows that flows are periodically analyzed accordingly to the flow timeout set in the controller. In this experiment, we set the flow timeout within 60 s. Therefore, the analysis machine receives the packets of each client every minute.

Two crucial aspects of this proposal are the time required to characterize the flow and the ability to increase the number of flows to be analyzed. As the machine only needs five packets to characterize flows, there is no need to wait until the end of the flow or connection to send the information to the threat detection architecture, which results in shorter detection times and, thus, a faster threat block. Furthermore, the analysis machine does not process all packets and then is immune to flooding attacks. Therefore, for this reason, the analysis machine can receive a much larger number of flows, ensuring robustness and potential for scalability.

The second experiment shows one of the client machines attacking the server. The attacker machine performs a threat that also spoofs the source IP to avoid detection. Similar to the first experiment, Figure 12 shows the traffic received by the server and the analysis machine. Our system immediately installs on all switches a rule blocking the source IP after the attack started, and the analysis machine detected it, sending the information to the detection application. This procedure blocks the malicious traffic received by the server, as shown in Figure 12(A). Nevertheless, when the attacker changes the IP for the first time, the flow is once again analyzed and blocked, as shown in Figure 12(B). Nonetheless, this time, when the controller receives the alarm, it starts to mark the path of all flows, so it can map in which port each flow enters the network. The system generates a new alert once more after the analysis, around 28 s because the attacker changes the source IP. When the controller receives this second alert, it blocks both the IP source and the port in which this flow enters the network. Here we choose to block the traffic from the TCP port in which the attacker enters the network. However, other policies are applicable, such as redirecting the traffic to filter out the malicious flow. From this point on, when the attacker changes once more the source IP address, the traffic is not even analyzed again since it is blocked in as near to its origin as possible.

7 | CONCLUSION

This article proposes a fast and accurate threat detection and prevention architecture. Machine learning algorithms perform threat detection using combined with stream processing. We implement five algorithms to detect known and unknown attacks and obtain high accuracy in classification, higher than 95%, and an excellent trade-off between threat detection and false-positive rates in anomaly detection. The results show that our architecture handles well both known attacks and detects zero-day attacks even without any previous information, only based on normal network behavior. Another contribution of this work is the creation of a security dataset to evaluate our proposal. Moreover, we also use another dataset with real data collected from broadband users of one of the most important network operators in Brazil. Both datasets contain real traffic data instead of simulation results. We measure the processing throughput of the proposed threat detection architecture, obtaining a threat detection time of about 4 μ s per sample, enabling prompt counter-measures against attackers and showing the scalability of our architecture. The proposed system achieves quick threat detection time thanks to the stream processing technology combined with machine learning algorithms in the lambda architecture.

Besides, we propose a scheme based on software defined networking and the periodic analysis of the first five packets of each flow in our threat *prevention* architecture. The scheme ensures a quick detection since it does not have to wait until the end of a flow to characterize it. The threat prevention architecture mirrors the traffic to sensors spread around the network. Therefore, the architecture does not add any delay to

user communication. Moreover, the scheme performs an effective threat block, even in scenarios in which the attacker uses spoofed IP addresses. The controller receives alerts from the threat detection application and installs rules to block the attacker's source IP. However, when an attacker changes its IP, the architecture detects it, based on the time difference between the alerts and maps the source of the attack to effectively block the threat. The proposed scheme is robust since both the controller and the analysis machine are protected against flooding attacks since not all packets are mirrored to the analysis machine and thus prevent overloading the controller. The proposed architecture scales well since its threat detection time is low, due to the stream processing core that improves its throughput when the parallelism increases. Furthermore, the schema that analyzes only a few packets of each flow prevents an attacker to flood sensor elements.

ACKNOWLEDGMENTS

This research is supported by CNPq, CAPES, FAPERJ, FAPESP (2014/50937-1, 2015/24485-9, 2018/23292-0), and NSF (DMS-2023495). The authors would like to thank Renato Souza Silva for the contribution to the Network Operator Dataset.

DATA AVAILABILITY STATEMENT

Data available on request from the authors.

ORCID

Antonio G. Pastana Lobato  <https://orcid.org/0000-0002-1544-2333>

Martin Andreoni Lopez  <https://orcid.org/0000-0002-4170-4341>

Alvaro A. Cardenas  <https://orcid.org/0000-0002-5142-9750>

Otto Carlos M. B. Duarte  <https://orcid.org/0000-0002-6642-4100>

Guy Pujolle  <https://orcid.org/0000-0003-4147-7270>

REFERENCES

1. Suthaharan S. Big data classification: Problems and challenges in network intrusion prediction with machine learning. *ACM SIGMETRICS Perform Eval Rev.* 2014;41(4):70-73.
2. Cardenas AA, Manadhata PK, Rajan SP. Big data analytics for security. *IEEE Secur Priv.* 2013;11(6):74-76.
3. Fernandes NC, Duarte OCMB. XNetMon: a network monitor for securing virtual networks. *IEEE International Conference on Communications (ICC).* IEEE; 2011:1-5.
4. Thangavel S, Kannan S. Detection and trace back of low and high volume of distributed denial-of-service attack based on statistical measures. *Concurr Comput Pract Exper.* 2019;1-22. <https://doi.org/10.1002/cpe.5428>
5. Kaspersky. *Kaspersky DDoS Intelligence Report for Q1 2016.* Kaspersky Lab.; 2016.
6. Fenil E, Mohan Kumar P. Survey on DDoS defense mechanisms. *Concurr Comput Pract Exper* 2020;32(4):e5114. <https://doi.org/10.1002/cpe.5114>.
7. Campista MEM, Rubinstein MG, Moraes IM, Costa LHMK, Duarte OCMB. Challenges and research directions for the future internet networking. *IEEE Commun Surv Tutor.* 2014;16(2):1050-1079.
8. Clay P. A modern threat response framework. *Netw Secur.* 2015;1(4):5-10.
9. Moreira MDD, Laufer RP, Fernandes NC, OCMB D. A stateless traceback technique for identifying the origin of attacks from a single packet. *IEEE International Conference on Communications (ICC).* IEEE; 2011:1-6.
10. Laufer RP, Velloso PB, de Oliveira CD, Moraes IM, Bicudo MDD, Duarte OCMB. Towards stateless single-packet IP traceback. *IEEE Conference on Local Computer Networks LCN'2007.* IEEE; 2007:548-555.
11. Buczak A, Guven E. A survey of data mining and machine learning methods for cyber security intrusion detection. *IEEE Commun Surv Tutor.* 2015;18(2):1153-1176.
12. Ji SY, Jeong BK, Choi S, Jeong DH. A multi-level intrusion detection method for abnormal network behaviors. *J Netw Comput Appl.* 2016;62:9-17.
13. Lakhina A, Crovella M, Diot C. Mining anomalies using traffic feature distributions. *SIGCOMM Comput Commun Rev.* 2005;35(4):217-228.
14. Fernandes G, Carvalho LF, Rodrigues JJPC, Proença ML. Network anomaly detection using IP flows with principal component analysis and ant colony optimization. *J Netw Comput Appl.* 2016;64:1-11.
15. Palmieri F, Fiore U, Castiglione A. A distributed approach to network anomaly detection based on independent component analysis. *Concurr Comput Pract Exper.* 2014;26(5):1113-1129.
16. Ringberg H, Soule A, Rexford J, Diot C. Sensitivity of PCA for traffic anomaly detection. Paper presented at: ACM SIGMETRICS, San Diego, CA, USA; 2007:109-120.
17. Amaral AA, de Souza Mendes L, Zarpelão BB, Junior MLP. Deep IP flow inspection to detect beyond network anomalies. *Comput Commun.* 2017;98:80-96.
18. Chellammal P, Malarchelvi SKPD. Real-time anomaly detection using parallelized intrusion detection architecture for streaming data. *Concurr Comput Pract Exper.* 2020;32:e5013.
19. Villar-Rodríguez E, Del Ser J, Torre-Bastida AI, Bilbao MN, Salcedo-Sanz S. A novel machine learning approach to the detection of identity theft in social networks based on emulated attack instances and support vector machines. *Concurr Comput Pract Exper.* 2016;28(4):1385-1395. <https://doi.org/10.1002/cpe.3633>
20. Li B, Zhang S, Li K. Towards a multi-layers anomaly detection framework for analyzing network traffic. *Concurr Comput Pract Exper.* 2017;29(14):e3955.
21. Bostani H, Sheikhan M. Hybrid of anomaly-based and specification-based IDS for Internet of Things using unsupervised OPF based on MapReduce approach. *Comput Commun.* 2017;98:52-71.
22. Singh K, Guntuku SC, Thakur A, Hota C. Big data analytics framework for peer-to-peer botnet detection using random forests. *Inf Sci.* 2014;278:488-497.

23. Viegas E, Santin A, Bessani A, Neves N. BigFlow: real-time and reliable anomaly-based intrusion detection for high-speed networks. *Future Gener Comput Syst.* 2019;93:473-485.
24. Kamburugamuve S, Govindarajan K, Wickramasinghe P, Abeykoon V, Fox G. Twister2: design of a big data toolkit. *Concurr Comput Pract Exper.* 2020;32(3):e5189. <https://doi.org/10.1002/cpe.5189>
25. Abeykoon V, Kamburugamuve S, Govindrarajan K, et al. Streaming machine learning algorithms with big data systems. *IEEE International Conference on Big Data.* IEEE; 2019:5661-5666.
26. Andreoni Lopez M, Mattos DMF, Duarte OCMB, Pujolle G. Toward a monitoring and threat detection system based on stream processing as a virtual network function for big data. *Concurr Comput Pract Exper.* 2019;31(20):e5344.
27. Cardenas AA, Baras JS & Seamon K A framework for the evaluation of intrusion detection systems. Paper presented at: IEEE Symposium on Security and Privacy (SP'06), Oakland, California, USA; 2006:15-77.
28. Lee W, Stolfo SJ, Mok KW. Mining in a data-flow environment: experience in network intrusion detection. *Proceedings of the Fifth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining.* ACM; 1999:114-124.
29. Sangkatsanee P, Wattanapongsakorn N, Charnsripinyo C. Practical real-time intrusion detection using machine learning approaches. *Comput Commun.* 2011;34(18):2227-2235.
30. Amini M, Jalili R, Shahriari HR. RT-UNNID: a practical solution to real-time network-based intrusion detection using unsupervised neural networks. *Comput Secur.* 2006;25(6):459-468.
31. Andreoni Lopez M, Mattos DMF, OCMB D. An elastic intrusion detection system for software networks. *Annals of Telecommunications.* 2016;71(11):595-605.
32. Andreoni Lopez M, Duarte OCMB. Providing elasticity to intrusion detection systems in virtualized software defined networks. Paper presented at: IEEE International Conference on Communications (ICC'15), London, England, UK; 2015:7120-7125.
33. Andreoni Lopez M. *A Monitoring and Threat Detection System Using Stream Processing as a Virtual Function for Big Data.* PhD thesis. Universidade Fedral do Rio de Janeiro and Sorbonne Université; 2018.
34. Lin YD, Lin PC, Yeh CH, Wang YC, Lai YC. An extended SDN architecture for network function virtualization with a case study on intrusion prevention. *IEEE Netw.* 2015;29(3):48-53.
35. Sharma PK, Singh S, Park JH. OpCloudSec: open cloud software defined wireless network security for the Internet of Things. *Comput Commun.* 2018;122:1-8.
36. Vicentini C, Santin A, Viegas E, Abreu V. SDN-based and multitenant-aware resource provisioning mechanism for cloud-based big data streaming. *J Netw Comput Appl.* 2019;126:133-149.
37. Lobato AGP, Andreoni Lopez M, Sanz IJ, Cardenas AA, Duarte OCMB, Pujolle G. An adaptive real-time architecture for zero-day threat detection. *IEEE International Conference on Communications (ICC).* IEEE; 2018:1-6.
38. Rychly M, Koda P & Smrz P Scheduling decisions in stream processing on heterogeneous clusters. Paper presented at: Eighth International Conference on Complex, Intelligent and Software Intensive Systems (CISIS), Birmingham City University, Birmingham, UK; 2014:614-619.
39. Marz N, Warren J. *Big Data: Principles and Best Practices of Scalable Realtime Data Systems.* 1st ed. Manning Publications Co.; 2013.
40. Andreoni Lopez M, Lobato A, Duarte OCMB. A performance comparison of open-source stream processing platforms. Paper presented at: IEEE GLOBECOM, Washington, USA; 2016:1-6.
41. Cheng Z, Caverlee J, Lee K. You are where you tweet: a content-based approach to geo-locating Twitter users. *Proceedings of the 19th ACM International Conference on Information and Knowledge Management CIKM'10.* ACM; 2010:759-768.
42. Lippmann RP, Fried DJ, Graf I, et al. Evaluating intrusion detection systems: the 1998 DARPA off-line intrusion detection evaluation. *Proceedings of DARPA Information Survivability Conference and Exposition. DISCEX'00.* Vol 2. IEEE; 2000:12-26.
43. Tavallaee M, Bagheri E, Lu W, Ghorbani AA. A detailed analysis of the KDD CUP 99 data set. *Proceedings of the Second IEEE Symposium on Computational Intelligence for Security and Defence Applications.* IEEE; 2009:1-6.
44. Sommer R, Paxson V. Outside the closed world: On using machine learning for network intrusion detection. *IEEE Symposium on Security and Privacy (SP).* IEEE; 2010:305-316.
45. Garcia S, Grill M, Stiborek J, Zunino A. An empirical comparison of botnet detection methods. *Comput Secur.* 2014;45:100-123.
46. Liu Q, Li P, Zhao W, Cai W, Yu S, Leung VC. A survey on security threats and defensive techniques of machine learning: a data driven view. *IEEE Access.* 2018;6:12103-12117.
47. Moraes IM, Mattos DMF, Ferraz LHG, et al. FITS: a flexible virtual network testbed architecture. *Comput Netw.* 2014;63:221-237. Special Issue on Future Internet Testbeds Part {II}.

How to cite this article: Lobato AGP, Andreoni Lopez M, Cardenas AA, Duarte OCMB, Pujolle G. A fast and accurate threat detection and prevention architecture using stream processing. *Concurrency Computat Pract Exper.* 2021;e6561. <https://doi.org/10.1002/cpe.6561>