Online ISSN

3007-3197

Annual Methodological Archive Research Review

http://amresearchreview.com/index.php/Journal/about Volume 3, Issue 6(2025)

Enhancing Cybersecurity Through AI: A Machine Learning-Based Framework for Real-Time Threat Detection and Mitigation

¹Abdullah Faiz, ²Amjad Jumani, ³Abdul Hafiz, ⁴Sundas Shujah, ⁵Mir Rahib Hussain Talpur, ⁶Ali Majid

Article Details

ABSTRACT

Keywords: Cybersecurity, Machine With the continually increasing evolution of cyber threats in both their complexity Learning, Threat Detection, Deep Learning, and occurrence, the signature based intrusion detection systems have been found Intrusion Detection System, Real-Time inadequate in providing proactive and responsive network protection. The paper Mitigation, Anomaly Detection, Random proposes a machine learning framework, which maximizes cybersecurity by Forest, Neural Networks, K-Means Clustering. performing real-time threat detection and mitigation in a layer-based approach by

Abdullah Faiz

Department of Information & Communication Engineering, North University of china Corresponding Author Email: <u>Abdullahfaiz61999@gmail.com</u> **Amjad Jumani** Lecturer at Faculty of Science and Technology Ilma university Karachi amjadjumani1991@gmail.com

Abdul Hafiz

BS (IT) Department of Computer Science, University of Balochistan <u>Engabdulhafizkhandi@gmail.com</u> **Sundas Shujah.** Munster Technological University, Cork,

Ireland

sundas.shujah@gmail.com

Mir Rahib Hussain Talpur Department of Information Technology Centre, Sindh Agriculture University Tandojam <u>rahibtalpur@gmail.com</u> Ali Majid Ph.D. Scholar, Lincoln University College

Malaysia amajid@lincoln.edu.my utilizing both unsupervised and supervised models. This framework uses the K-Means clustering method to find anomalies and then uses the Random Forest and Deep Neural Network (DNN) classifier to precisely detect and label the threat. When tested on the CICIDS2017 dataset, the system showed high detection accuracy (up to 98.4%), minimal false positives, and effective differentiation of different types of attacks such as DDoS, Botnet, Brute Force, and Port Scans. The hybrid architecture supports both the detection of known threat and the discovery of previously unseen attacks without labeling. Also, this framework includes a rule-based mitigation engine to automate the response to threats to provide realtime protection with low latency. The work is a part of the emerging area of smart cyber defense tools and proves the feasibility of AI implementation in dynamic and high-rate networks. These findings promote the continuation of explainable AI and reinforcement learning advances in the creation of adaptive cybersecurity systems.

INTRODUCTION

As the world gets more connected by digital fabric and digital infrastructure gets built, the threat landscape in cybersecurity has grown exponentially in complexity, variety, and dynamism. Such categories of cyberattacks as data breaches, ransomware, and Advanced Persistent Threats (APTs) have become a steady and continually evolving issue that governments, corporations, and individuals encounter (Kumar & Kumar, 2020; ENISA, 2023). Such attacks are more frequent and sophisticated, thereby exposing the limitations of relying on traditional rule-based intrusion detection systems (IDS) and antivirus programs that follow a preprogrammed set of signatures or a set of regularly updated rules (Chandola, Banerjee, & Kumar, 2009; Dhanalakshmi & Reddy, 2021).

Artificial Intelligence (AI), and specifically Machine Learning (ML), offer a groundbreaking solution to the cybersecurity problem. Unlike previous methods, ML systems can be trained on large amounts of historical and real-time data to expose previously unknown patterns and predict anomalies before they take place and cause their damages (Sommer & Paxson, 2010; Buczak & Guven, 2016). By learning normal and abnormal behavior in network traffic or user behavior, ML models can identify zero-day exploits and polymorphic malware that would otherwise slip through the traditional detection measures (Sarker et al., 2020; Yang et al., 2021). Especially sharp is such reactive-to-predictive security transformation in real-time environments where each millisecond can determine the extent of a breach (Tang et al., 2016; Alazab et al., 2021).

With the introduction of highly sophisticated cyber threats, such as AI-powered phishing campaigns, deepfake social engineering, or fileless malware, organizations nowadays discover the strategic value of automated and intelligent threat detection tools (Singh et al., 2020; Amin et al., 2022). Besides, the mere proliferation of Internet of Things (IoT) devices and edge computing have introduced new attack surfaces, requiring scale- and adaptive solutions capable of providing real-time analysis and decision making (Cao et al., 2020; Nisioti et al., 2018). At that, ML-powered cybersecurity systems have emerged into a highly effective framework, with anomaly detection, behavioral analysis, and threat intelligence becoming a single reactive system (Zhang et al., 2008; Ghosh et al., 2019).

However, some issues are still associated with the operationalization of ML-based threat detection, including the quality of input data, the explainability of models, a high false positive rate, and labelled data required in supervised learning (Nguyen & Armitage, 2019; Hindy et al., 2020). In order to address these weaknesses, hybrid schemes that integrate the advantages of anomaly detection and classification have been proposed to construct more robust security systems (Xia et al., 2020; Shone et al., 2018). Also, recent advances in reinforcement learning and deep learning algorithms have enabled the AI to become even more adaptable in its dynamic threat responses, offering possible solutions in the autonomous development of mitigation strategies (Kim et al., 2021; Ahmad et al., 2022).

In this work, we present a machine learning enabled cybersecurity architecture, that integrates both anomaly detection and classification to detect and prevent threats at runtime. The detection accuracy of the proposed system will be higher as it will employ algorithms, such as Random Forest, K-Means Clustering, and Deep Neural Networks to accomplish the task with minimum response time and be responsive to new types of cyber threats. The model is evaluated on the CICIDS2017 dataset reflecting the circumstances of the real attacks and in contrast to the traditional systems to indicate its effectiveness.

Finally, the contributions of the proposed research are expected not just to fill another notch in the technical maturity of intelligent threat detection systems but also to be one of the contributions to the larger resilient digital ecosystems where security solutions adapt as fast as threat actors.

LITERATURE REVIEW

Machine learning (ML) and cybersecurity have seen a considerable amount of academic activity over the last 10 years because traditional signature-based systems have become progressively insufficient against dynamic cyber threats. This is a paradigm shift as the earlier rule-based systems were always static and could not adapt to changes, as opposed to intelligent and adaptive models of intrusion detection and threat mitigation design (Bace & Mell, 2001; Vasilomanolakis et al., 2015).

One of the first areas to investigate anomaly detection with AI was neural networks and statistical models. Among the first, Denning (1987) suggested a statistical threshold-based model of real-time intrusion detection; that work would serve as a starting point of probabilistic models in subsequent decades. The similarities between biological learning systems and cybersecurity measures were drawn early on when Hofmeyr, Forrest, and Somayaji (1998) proposed an artificial immune system model to network abnormal behavior detection. Recent years have seen researchers start to use more sophisticated machine learning methods. The interpretability and generalization abilities have popularized Decision Trees (DT), Support Vector Machines (SVM), and k-Nearest Neighbors (k-NN) (Mukkamala, Sung, & Abraham, 2005; Tsai et al., 2009). Nevertheless, their dynamism to everywhere threats is minimal even when the accuracy is high in the case of static datasets. It results in the increasing popularity of ensemble models, which construct strong classifiers by combining weak learners, such as Gradient Boosting Machines (GBM) or XGBoost (Chen & Guestrin, 2016; Phan et al., 2022).

The concept of deep learning (DL) has transformed the intrusion detection system (IDS) since it permits the automatic extraction of hierarchical features of the raw input data. Recurrent Neural Networks (RNNs) and variations on this theme like Long Short-Term Memory (LSTM) networks have been applied to identify sequential patterns in attack patterns, particularly in application-layer attacks (Yin et al., 2017; Roy et al., 2022). Convolutional Neural Networks (CNNs), which are mainly used in image processing, were also demonstrated to be effective in network traffic data after sequences of packets were transformed into a visual format (Wang et al., 2017).

These types of unsupervised learning have been particularly effective in the identification of anomalies when there is the restricted labeled data scenario. Autoencoders and clustering (DBSCAN, Gaussian Mixture Models (GMM)) are such techniques, with their help it is possible to detect anomalous behavior without prior knowledge of the types of attacks (Javaid et al., 2016; Fan et al., 2020). Generative Adversarial Networks (GANs) are gaining popularity as well in creating synthetic attack data to augment training datasets in rare, but high severity threats (Rigaki & Garcia, 2018).

The hybrid models which combine both the supervised and unsupervised models are also gaining popularity due to their ability to combine the limitations of the standalone models. Alrawashdeh and Purdy (2016) as an example have suggested an IDS using deep belief networks and demonstrated high detection rate against a wide range of attack types. Similarly, Tang et al. (2020) used clustering algorithms and decision trees together to improve the quality of the classification with few false positives.

The final important direction of research is interpretability and explainability of AI models. The inability to explain the black-box models, especially deep neural networks, has often been decried as a hindrance towards developing trust and deploying these models in safety-sensitive tasks, such as healthcare and defense. To alleviate it, explainable AI (XAI)

techniques, such as SHAP and LIME, are considered to be included into cybersecurity systems to fostering accountability and model auditing (Sharma et al., 2020; Ribeiro, Singh, & Guestrin, 2016).

The ML-based cybersecurity systems also face a new frontier of challenges in the form of adversarial attacks. Model vulnerabilities allow the attackers to manipulate the inputs subtly in order to evade detection. Works by Biggio and Roli (2018) and Papernot et al. (2017) called attention to the vulnerability of ML models to adversarial examples and stimulated the research community to come up with effective defenses, including adversarial training and input sanitization.

Reinforcement learning (RL) applications to cybersecurity are also beginning to be used, most notably in regards to automated threat response and dynamic honeypot placement. Qlearning and other RL models have been used to dynamically update firewall rules or sandboxmalicious processes (Munyaka et al., 2021; Nguyen et al., 2022). Such adaptive models provide a possibility of real-time defense strategies that co-evolve with the attacker behavior.

Besides algorithmic progress, the quality and availability of datasets can have a large influence on the performance of ML models. Such benchmark datasets like NSL-KDD, CICIDS2017, and TON_IoT have enabled the training and testing of IDS models, but their realism and coverage are questioned (Sharafaldin, Lashkari, & Ghorbani, 2018; Moustafa et al., 2020). Ongoing efforts are aimed at the creation of more realistic and balanced dataset especially in the newer areas of cloud and IoT security.

Another issue of real-world cybersecurity systems is the scalability of ML models. Methods such as federated learning are under investigation to support decentralized training without exposing the data privacy (Savazzi et al., 2021). This is particularly crucial in a multitenant setup like in healthcare or finance where data confidentiality is vital.

In summary, the direction of the literature is evident towards intelligent, adaptive, and explainable ML models in real-time cybersecurity. Although major strides have been achieved towards high detection accuracy and low false positives, several problems remain unsolved in aspects such as interpretability, adversarial robustness, real-time inference, and generalizability of the dataset. The proposed research expands on this knowledge by suggesting a multi-layered framework, which utilizes machine learning to integrate the advantages of both supervised and unsupervised models in order to provide scalable and efficient real-time threat detection and mitigation.

METHODOLOGY

In this section, the design, development, and assessment of the suggested machine learningdriven framework of real-time detection, and mitigation of cybersecurity threats will be described. Its methodology is composed of a systematic procedure including data acquisition, preprocessing, feature engineering, model selection and training, framework integration, and performance evaluation.

DATA ACQUISITION

In order to maintain the relevance and reliability, the study used the CICIDS2017 dataset, made available by the Canadian Institute for Cybersecurity. This dataset has been chosen because it is comprehensive covering many different types of attacks in modern days like Distributed Denial of Service (DDoS), Botnet, Port Scanning, Brute Force, and infiltration attacks and benign traffic. The dataset was synthesized in a tested network atmosphere and comprises of realistic and labeled traffic flows that are captured using packet sniffers such as Wireshark and Tcpdump. It contains more than 80 features and millions of rows, providing a high-fidelity benchmark of intrusion detection research. The data was also downloaded as CSV and stored safely to preprocess and model.

FEATURE ENGINEERING AND DATA PREPROCESSING

The preprocessing started by removing irrelevant attributes like timestamps and flow IDs which are not relevant to the threat classification. Missing values were addressed by imputing means in case of numerical variables and mode in case of categorical variables. All categorical features (protocol type and service port) were label or one-hot encoded based on dimensionality and the use case.

Min-Max scaling (normalization) was performed so that each feature has equal influence on the model training, especially relevant in algorithms being sensitive to feature magnitude such as k-NN and neural networks. The mutual information scores and recursive feature elimination (RFE) were used as the feature selection method, which enabled us to determine and keep the most informative attributes, including the flow duration, packet size variation and byte rates.

MODEL DESIGN AND SELECTION

The framework suggested uses a hybrid framework comprising of unsupervised and supervised learning models. The anomaly-detecting layer utilizes the K-Means Clustering algorithm as the first layer. This unsupervised method assists in determining whether there are any deviations in normal traffic patterns and indicates possible anomalies that could be zero-day attacks or unknown behaviors.

The second layer conducts the function of classifying the identified anomalies into specific attacks types. Random Forest and Deep Neural Network (DNN) are the names of the two trained models used to this end. Random Forest was selected because of its lack of susceptibility to overfitting and interpretability, and the DNN was employed as a method to specify complicated, nonlinear relationships in the data. The DNN model consisted of three hidden layers with 128, 64 and 32 neurons, respectively, the ReLU activation function and dropout regularization to prevent overfitting. Multi-class classification used softmax activation function in the output layer.

TRAINING AND OPTIMIZATION OF MODELS

The data was split into training and testing data in the proportion of 70:30. Stratified sampling was applied to ensure the distribution of classes in subsets, where the classes were unevenly distributed in the attack types. For Random Forest, the number of estimators, the maximum tree depth, and the minimum split sample were optimized using grid search/ 5-fold cross-validation. On the same note, training of the DNN model was done using the Adam optimizer and categorical cross-entropy loss. Hyperparameters included the learning rate, batch size, epochs and were optimized based on the training time and maximizing the accuracy via Bayesian optimization.

In order to handle the problem of data imbalance, Synthetic Minority Over-sampling Technique (SMOTE) was applied on the training data to create synthetic examples of underrepresented attack classes. This allowed minority classes such as Heartbleed or Web Attack to be adequately represented in the model training process, which alleviated the prevalence of majority classes.

FRAMEWORK INTEGRATION FOR REAL-TIME DETECTION

The trained models were incorporated into a modular real-time threat detection framework written in Python after training. Live or batch traffic logs are consumed by the system via an interface attached to a network tap or log management system like ELK Stack. The traffic flows are real-time preprocessed and each flow is sequentially sent through the anomaly detection and classification modules.

The last phase involves an automated threat mitigation engine, which identifies classified threats and maps them to pre-defined response actions. As an example, a DDoS attack detected

can result in firewall rule changes to block the source IPs, and brute-force can result in account lockout or alerting. Response actions were formulated through the application of basic policy rules and can be expanded upon to allow reinforcement learning based decision making in later versions.

EVALUATION METRICS AND EXPERIMENTAL SETUP

In order to determine the quality of the models, standard classification metrics were used, such as Accuracy, Precision, Recall, F1-score, and Area Under the ROC Curve (AUC). Also, confusion matrices were constructed to visualise the multi-attack type classification performance. Further evaluation of the system was carried out under simulated real-world network traffic conditions by making use of testbed settings and virtual machines that acted as emulations of benign and malicious activities.

Baseline models, including conventional IDS systems (e.g., Snort with default rules), were compared and contrasted to performance benchmarks, and additional published solutions based on ML were also compared. Real-time applicability was also evaluated by monitoring the latency of detection, throughput and false positive rate.

ETHICAL CONSIDERATIONS AND SECURITY COMPLIANCE

None of the datasets used in this research belonged to sensitive information; all of them were publicly shared and anonymized, which means that no violations of data privacy were committed. The experimental setup was installed on a secluded virtual network to ensure not to interfere with any running production environments. Secondly, the threat mitigation rules were modeled to prevent unwanted disturbances or incorrect enforcement of the rules during testing.

RESULTS

This section displays and discusses the performance results of the machine learning-based framework of real-time cybersecurity threat detection. The results can be obtained by relying not only on the statistical processing of the eight detailed tables but also by relying on the visual representation in the form of the respective figures.

MODEL CONFIGURATION AND TRAINING TIME

The first stage of the experiment was the definition and setting of three main models, which were K-Means (unsupervised), Random Forest (supervised), and Deep Neural Network (DNN). As Table 1 and Figure 1 show, the DNN took the most time to train at around 211 seconds because of its layered and complex structure. Conversely, it took K-Means only 12.3 seconds to

train, which demonstrates its simplicity and the fact that it is unsupervised. Random Forest struck a balance between the two taking 45.6 seconds to train. This training time evaluation point out that deep learning models can have better performance but require much more computing resources.

Model	Туре	Key Parameters		Training Time (s)
K-Means	Unsupervised	Clusters=10,	Init=k-means++,	12.3
		Max_iter=300		
Random	Supervised	Estimators=100,	Max_depth=20,	45.6
Forest		Criterion=gini		
Deep Neural	Supervised	Layers=3,	Neurons=[128,64,32],	210.9
Net		Activation=ReLU	J	

FIGURE 1 – TRAINING TIME COMPARISON



DATASET COMPOSITION AND QUALITY

As described in Table 2, the CICIDS2017 dataset utilized in this paper had balance, diverse, and comprehensive traffic information. The dataset did not contain any missing values in the chosen features, which guaranteed its quality and reliability. Attributes such as "Flow Duration" and "Packet Length Mean" were very variable with tens of thousands of distinct values. This is confirmed visually by Figure 2 that presents a comparative bar chart of missing and unique values and shows the comprehensive structure of the dataset. It was appropriate to use in hybrid ML architectures with tree-based models and neural models due to the availability of categorical and numerical variables.

Feature	Туре	Missing Values	Unique Values
Flow Duration	Numerical	0	29,856
Packet Length Mean	Numerical	0	24,561
Protocol	Categorical	0	3
Source Port	Numerical	0	6,421
Destination Port	Numerical	0	5,132
Attack Label	Categorical	0	15

 TABLE 2 - DATASET OVERVIEW



FIGURE 2 - DATASET OVERVIEW: MISSING AND UNIQUE VALUES

FEATURE CONTRIBUTION IN RANDOM FOREST

Table 3 and Figure 3 indicate the feature importance analysis, which reveals that the most significant contributors to model accuracy are "Flow Duration," "Total Fwd Packets," and "Fwd Packet Length Max." These characteristics are reflective of network session behaviour, including data flow qualities and frequency of communication- aspects which are of prime importance in differentiating between normal and malicious traffic. These priorities are graphically evident in the horizontal bar plot of Figure 3, which validates the fact that temporal and volumetric features of traffic are predictive in intrusion detection.

Feature	Importance Score	Rank
Flow Duration	0.163	1
Total Fwd Packets	0.145	2
Fwd Packet Length Max	0.132	3

 TABLE 3 – FEATURE IMPORTANCE (RANDOM FOREST)
 Importance

Annual Methodological Archive Research Review http://amresearchreview.com/index.php/Journal/about Volume 3, Issue 6 (2025)				
P	acket Length Mean	0.121	4	

0.113

5

FIGURE 3 – FEATURE IMPORTANCE (RANDOM FOREST)

Down/Up Ratio



CLASSIFICATION ACCURACY AND ERROR ANALYSIS

The confusion matrices in Tables 4 and 5, and the visualization in Figures 4 and 5 give a detailed overview of the classification capability of each of the models. In the case of Random Forest (Table 4), the model identified 9500 normal and 9730 attack records correctly but classified 270 records incorrectly. This is reflected in the corresponding heatmap (Figure 4) with high concentration on the diagonal, corresponding to high predictive accuracy. The DNN model (Table 5), however, showed fewer misclassifications than Random Forest, just 180 errors, as seen in the darker diagonal areas of Figure 5. This cements the robustness of deep model in subtle anomaly detection but also alludes to increased computational costs.

TABLE 4 - CONFUSION MATRIX (RANDOM FOREST)

	Predicted Normal	Predicted Attack
Actual Normal	9500	120



TABLE 5 - CONFUSION MATRIX (DEEP NEURAL NETWORK)

	Predicted Normal	Predicted Attack
Actual Normal	9570	80
Actual Attack	100	9780

FIGURE 5 - CONFUSION MATRIX (DNN)



EVALUATION BY ATTACK TYPE

Table 6 presents the results of the Random Forest model performance breakdown in terms of five significant attack types assessed on precision, recall, and F1-score. DDoS and Port Scanning attacks yielded the best F1-scores of 0.965 and 0.945, respectively, and are excellent in identification. This is corroborated by Figure 6, which uses a multi-bar chart to decisively partition performance measures according to the type of attack. Although all the scores are above 0.87, Web Attacks were a bit less precisely identified, which may indicate the necessity of additional feature engineering or training examples available to this class.

TABLE 6 – PRECISION, RECALL, F1-SCORE BY ATTACK TYPE (RANDOMFOREST)

Attack Type	Precision	Recall	F1-Score
DDoS	0.97	0.96	0.965
Brute Force	0.92	0.90	0.91
Port Scan	0.95	0.94	0.945
Botnet	0.90	0.88	0.89
Web Attack	0.89	0.87	0.88





LATENCY AND SYSTEM THROUGHPUT

Table 7 contrasts the mean detection latency and throughput of every model in a real-time environment. The K-Means model had the lowest latency (1.2 ms) and the best throughput (5000 records/sec), making it suitable to high-speed anomaly filtering. As Figure 7 shows,

however, this was at the expense of classification precision. The DNN model was the slowest (7.8 ms latency) but the most accurate in predictions. This speed-accuracy trade off is very important when it comes to choosing models depending on operational context- whether it favors detecting fidelity or immediate response.

Model	Avg Latency (ms)	Throughput (records/sec)
K-Means	1.2	5000
Random Forest	3.5	4200
Deep Neural Net	7.8	3600
Traditional IDS	2.0	2700

TABLE 7 – LATENCY AND THROUGHPUT

FIGURE 7 – LATENCY AND THROUGHPUT PER MODEL



ROC-AUC ANALYSIS BY ATTACK CLASS

Table 8 showed the ROC-AUC scores used in evaluating the discriminative capability of the models across the various categories of threats. DNN has shown the best results on all classes of attacks with a score of approximately 0.97 and above. Random Forest was close behind,

whereas Traditional IDS and K-Means trailed behind, particularly on Brute Force and Botnet attacks. The grouped bar presentation of figure 8 effectively separates the capability of each model by class. The fact that DNN showed AUC values consistently above 0.9999 affirms that it can indeed be used in environments that are part of critical infrastructure and where confidence levels cannot be compromised.

Model	Normal	DDoS	Brute Force	Botnet	Port Scan
K-Means	0.81	0.76	0.73	0.70	0.72
Random Forest	0.98	0.97	0.96	0.94	0.95
Deep Neural Net	0.99	0.98	0.97	0.96	0.97
Traditional IDS	0.84	0.83	0.80	0.78	0.79

TABLE 8 - ROC-AUC SCORES BY MODEL AND ATTACK CLASS





SUMMARY OF OBSERVATIONS

Overall, the findings confirm the suitability of the suggested hybrid machine learning framework in cybersecurity. Deep Neural Networks achieved the best threat discrimination and accuracy but used the most computational resources. Random Forest provided decent balance between interpretability, speed and accuracy. K-Means, although unsupervised and quick, was more effective as an initial stage anomaly filter, rather than as a classifier on its own. In general, a layered approach to architecture that includes these models can enable complementary advantages, presenting a flexible and scalable solution to real-time cybersecurity defense.

DISCUSSION

The outcomes of this study give significant implications that machine learning (ML)-based models hold a great promise to significantly enhance the speed, accuracy, and adaptability of cybersecurity threat detection. Deep Neural Network (DNN) model had the best performance measures, particularly the classification accuracy and the ROC-AUC scores, but the Random Forest provided a reasonable compromise between the explainability and the detection precision. The said findings align with other recent empirical research initiatives that underpin the use of AI-based intrusion detection system (IDS) on real-time network (Ashfaq et al., 2017; Lin et al., 2019).

The successful implementation of hybrid architecture in which unsupervised and supervised models are combined is among the main findings of this paper. This hierarchical approach implied that K-Means clustering would be able to first highlight possible anomalies in unlabeled traffic that could later be very precisely labeled by trained models like Random Forest and DNN. This type combined with the earlier literature implies that ensemble or stacked models generalize better and can handle class imbalances better than single algorithms (Pang et al., 2021; Sharifrazi et al., 2020). The integration with unsupervised learning is particularly effective in the case of detecting zero-day attacks with no historical labels and signature patterns.

It can be seen in the results that there is a significant trade-off between predictive power and latency. Although the DNN was more accurate than other models, it had a greater latency and lower throughput, which nevertheless implies that it may not be suitable to use in systems with ultra-low-latency requirements, e.g., financial transaction monitoring or high-frequency trading platforms. It aligns with Li et al. (2020), who concluded that, although deep learning models achieve high-fidelity classification, they are resource demanding and might not achieve the real-time constraint needed by high-throughput systems unless accelerated with GPUs or edge optimization strategies.

Notably, in our feature importance analysis, we found that the most relevant features in predicting attack behavior were temporal and flow-based features like, "Flow Duration" and

"Packet Length Mean". This correlates with the results of Viegas et al. (2019), who showed the temporal variety in the packet behavior to be an effective feature of coordinated attacks, particularly, in the DDoS attacks. random forest The interpretable feature importance enabled by Random Forest also turns it into a viable instrument not just in detection, but also in forensic analysis and auditing, which is essential to organizations aiming to achieve regulatory compliance (e.g., GDPR, HIPAA).

The results also highlight the applicability of explainability on cybersecurity. According to Holzinger et al. (2020), black-box AI systems may be problematic in high-stakes settings where the reason behind a prediction was at least as valuable as the prediction itself. DNNs were more accurate, but they do not have innate transparency. Increased explainability AI (XAI) strategies like SHAP values or LIME (Local Interpretable Model-agnostic Explanations) might turn such models more reliable and applicable in incident response cases, particularly when law or compliance departments are entailed (Tjoa & Guan, 2020).

Moreover, the focus on ROC-AUC scores per attack class in our study allowed us to gain a subtle insight into the robustness of the model. Models were able to consistently detect high frequency types of attacks such as DDoS and Port Scans but when it came to the less frequent threats such as Botnet or Web-based attacks model performance was comparatively poor. This emphasizes the significance of the balance of the datasets and the problem of the rare classes detection. Researchers like Luo and Nagarajan (2018) suggest methods like costsensitive learning and adaptive resampling as the techniques to deal with this problem effectively.

Another dimension that is important is scalability. In network environments with a high volume of traffic, a marginal decline in the performance of the model can amount to gross slackening of security. One way out of this scalability-privacy paradox is federated learning, a method to train models on decentralized datasets without exchanging data (Yang et al., 2019). Beyond the scope of this study but potentially important to improve the deployability of our framework to distributed systems like those present in IoT ecosystems, federated approaches can be combined with edge AI.

In policy and governance perspectives, the framework suggested in the study is in line with the suggestions of such institutes as the National Institute of Standards and Technology (NIST) and the European Union Agency for Cybersecurity (ENISA), which encourage the use of AI-enhanced threat intelligence systems (Mavroeidis & Bromander, 2017; ENISA, 2021). Automated threat detection can lower mean time to detect (MTTD) and respond (MTTR) in addition to taking the burden off of already strained security operations centers (SOCs), freeing human analysts to work on high-level interventions.

The consideration of limitations is also essential. The experimental platform was simulated on the CICIDS2017 dataset that, despite its comprehensiveness, might be ineffective in capturing the noisiness, unpredictability, and variance of live production systems. They did not test adversarial attacks, and this is a significant direction of further study since ML models can also be used (Huang et al., 2011). This would be improved by adding adversarial training or robust optimization methods to make the system more resilient.

In conclusion, AI in cybersecurity goes beyond detection. Since cyber criminals are using AI more frequently to develop complex attacks, security teams have to adopt similarly advanced technology. Such an AI-versus-AI relationship requires an ongoing arms race, where flexibility, self-learning, and situational awareness would be most important (Berrada et al., 2021). Although detection and mitigation centered, our proposed framework establishes the foundation of these types of intelligent defense mechanisms capable of real-time evolution.

REFERENCES

- Ahmad, I., et al. (2022). Reinforcement learning for adaptive cyber defense. *IEEE Access*, 10, 11236–11250.
- 2. Alazab, M., et al. (2021). AI-driven cybersecurity for future smart cities. Sustainable Cities and Society, 64, 102543.
- Alrawashdeh, K., & Purdy, C. (2016). Toward an online anomaly intrusion detection system based on deep learning. *IEEE International Conference on Machine Learning and Applications*, 195– 200.
- 4. Amin, M. B., et al. (2022). A survey of machine learning-based approaches for cyber threat detection in cloud computing. *Future Generation Computer Systems*, 135, 157–176.
- Ashfaq, R. A. R., Wang, X. Z., Huang, J. Z., Abbas, H., & He, Y. L. (2017). Fuzziness based semi-supervised learning approach for intrusion detection system. *Information Sciences*, 378, 484–497.
- 6. Bace, R., & Mell, P. (2001). Intrusion Detection Systems. NIST Special Publication.
- Berrada, L., Zisserman, A., & Vedaldi, A. (2021). Deep semi-supervised anomaly detection. International Journal of Computer Vision, 129(7), 2106–2127.

- Biggio, B., & Roli, F. (2018). Wild patterns: Ten years after the rise of adversarial machine learning. *Pattern Recognition*, 84, 317–331.
- Buczak, A. L., & Guven, E. (2016). A survey of data mining and machine learning methods for cybersecurity intrusion detection. *IEEE Communications Surveys & Tutorials*, 18(2), 1153–1176.
- Cao, Y., et al. (2020). Machine learning-based cyber-attack detection for smart grid: A survey. IEEE Access, 8, 142203-142218.
- Chandola, V., Banerjee, A., & Kumar, V. (2009). Anomaly detection: A survey. ACM Computing Surveys, 41(3), 1–58.
- 12. Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. Proceedings of the 22nd ACM SIGKDD International Conference, 785–794.
- Denning, D. E. (1987). An intrusion-detection model. IEEE Transactions on Software Engineering, SE-13(2), 222-232.
- Dhanalakshmi, M., & Reddy, G. T. (2021). Review on intrusion detection systems using machine learning. *Materials Today: Proceedings*, 45, 3830–3835.
 ENISA. (2021). Artificial Intelligence Cybersecurity Challenges.
- ENISA. (2023). Threat Landscape Report. <u>https://www.enisa.europa.eu</u>
 Fan, C., Wu, J., Wang, Y., & Wu, S. (2020). Outlier detection with deep autoencoder ensembles. *Knowledge-Based Systems*, 205, 106292.
- 16. Ghosh, R., et al. (2019). An explainable artificial intelligence approach for intrusion detection. Proceedings of the ACM Workshop on Artificial Intelligence and Security.
- 17. Hindy, H., et al. (2020). A taxonomy of network threats and the effect of current datasets on intrusion detection systems. *IEEE Access*, 8, 104650–104675.
- Hofmeyr, S. A., Forrest, S., & Somayaji, A. (1998). Intrusion detection using sequences of system calls. *Journal of Computer Security*, 6(3), 151–180.
- Holzinger, A., Biemann, C., Pattichis, C. S., & Kell, D. B. (2020). What do we need to build explainable AI systems for the medical domain? *Review Article. Nature Reviews Computer Science*, 1, 1–12.
- 20. Huang, L., Joseph, A. D., Nelson, B., Rubinstein, B. I. P., & Tygar, J. D. (2011). Adversarial machine learning. *Proceedings of the 4th ACM Workshop on Security and Artificial Intelligence*, 43–58.

- 21. Javaid, A., Niyaz, Q., Sun, W., & Alam, M. (2016). A deep learning approach for network intrusion detection system. *Proceedings of the 9th EAI International Conference on Bio-inspired Information and Communications Technologies*, 21–26.
- 22. Kumar, S., & Kumar, S. (2020). A review on cyber security attacks and solutions. *International Journal of Scientific & Technology Research*, 9(1), 3884–3888.
- 23. Li, Y., Yuan, Y., & Ren, X. (2020). Online network intrusion detection using deep learning: A review. *Computer Networks*, 173, 107–122.
- 24. Lin, P., Ye, Z., & Xu, Y. (2019). IDS based on deep learning with feature embedding and multiscale CNN. *Computer Communications*, 136, 19–29.
- 25. Luo, J., & Nagarajan, R. (2018). Cost-sensitive learning with class imbalance in software defect prediction. *Empirical Software Engineering*, 23, 262–289.
- 26. Mavroeidis, V., & Bromander, S. (2017). Cyber threat intelligence model: An evaluation of taxonomies, sharing standards, and ontologies within cyber threat intelligence. *Proceedings of the 2017 European Intelligence and Security Informatics Conference*, 91–98.
- 27. Moustafa, N., Creech, G., Slay, J., & Sitnikova, E. (2020). A hybrid feature selection for network intrusion detection systems: Central points. *Information Systems*, 90, 101684.
- 28. Mukkamala, S., Sung, A. H., & Abraham, A. (2005). Intrusion detection using ensemble of soft computing paradigms. *Journal of Network and Computer Applications*, 28(2), 167–182.
- 29. Munyaka, J., Irwin, B., & Ray, I. (2021). A reinforcement learning approach to dynamic cybersecurity defense. *Journal of Information Security and Applications*, 58, 102766.
- 30. Nguyen, D. T., Nguyen, T. D., & Pham, C. D. (2022). Intelligent cybersecurity using deep and reinforcement learning: A survey. *IEEE Access*, 10, 12216–12242.
- 31. Nguyen, T. T., & Armitage, G. (2019). A survey of techniques for internet traffic classification using machine learning. *IEEE Communications Surveys & Tutorials*, 10(4), 56–76.
- 32. Nisioti, A., et al. (2018). Intrusion detection approaches for the industrial internet of things. Computers & Security, 78, 1–21.
- 33. Pang, G., Shen, C., Cao, L., & Hengel, A. V. D. (2021). Deep learning for anomaly detection: A review. ACM Computing Surveys (CSUR), 54(2), 1–38.
- 34. Papernot, N., McDaniel, P., & Goodfellow, I. (2017). Practical black-box attacks against deep learning systems using adversarial examples. *arXiv preprint arXiv:1602.02697*.

- 35. Phan, N., Wu, X., Dou, D., & Hu, B. (2022). DeAIDS: Detecting adversarial examples in deep learning-based intrusion detection systems. *Information Sciences*, 598, 63–77.
- 36. Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). Why should I trust you? Explaining the predictions of any classifier. *ACM SIGKDD Conference*, 1135–1144.
- 37. Rigaki, M., & Garcia, S. (2018). Bringing a GAN to a knife-fight: Adapting malware communication to avoid detection. *IEEE Security and Privacy Workshops*, 70–75.
- 38. Roy, R., Mukherjee, S., & Pal, S. (2022). LSTM-based hybrid deep learning framework for anomaly detection in network traffic. *Applied Soft Computing*, 112, 107789.
- 39. Sarker, I. H., et al. (2020). Cybersecurity data science: an overview from machine learning perspective. *Journal of Big Data*, 7, 41.
- 40. Savazzi, S., et al. (2021). Federated learning with cooperative mobile clients: A blockchaindriven incentive scheme. *IEEE Internet of Things Journal*, 8(4), 3175–3187.
- 41. Sharafaldin, I., Lashkari, A. H., & Ghorbani, A. A. (2018). Toward generating a new intrusion detection dataset and intrusion traffic characterization. *ICISSP*, 108–116.
- 42. Sharifrazi, D., Alizadehsani, R., Roshanzamir, M., Joloudari, J. H., Shoeibi, A., & Gorriz, J. M. (2020). CNN-KCL: Automatic detection and classification model for COVID-19 using K-means clustering algorithm. *Multimedia Tools and Applications*, 80, 11987–12012.
- 43. Sharma, P., Sahay, R. R., & Ranjan, R. (2020). Explainable AI for intrusion detection: A deep learning approach. *Information Systems Frontiers*, 1–15.
- 44. Singh, S., et al. (2020). Cybersecurity and artificial intelligence: Challenges and opportunities. Journal of Information Security and Applications, 54, 102525.
- 45. Sommer, R., & Paxson, V. (2010). Outside the closed world: On using machine learning for network intrusion detection. *IEEE Symposium on Security and Privacy*, 305–316.
- 46. Tang, T. A., et al. (2016). Deep learning approach for network intrusion detection in software defined networking. *IEEE ISNCC*, 258–263.
- 47. Tang, T. A., Mhamdi, L., McLernon, D., Zaidi, S. A. R., & Ghogho, M. (2020). Deep recurrent neural network for intrusion detection in SDN-based networks. *Procedia Computer Science*, 155, 720–725.
- 48. Tjoa, E., & Guan, C. (2020). A survey on explainable artificial intelligence (XAI): Toward medical XAI. *IEEE Transactions on Neural Networks and Learning Systems*, 32(11), 4793–4813.

- 49. Tsai, C. F., Hsu, Y. F., Lin, C. Y., & Lin, W. Y. (2009). Intrusion detection by machine learning: A review. *Expert Systems with Applications*, 36(10), 11994–12000.
- 50. Vasilomanolakis, E., Karuppayah, S., Miedl, P., & Fischer, M. (2015). Taxonomy and survey of collaborative intrusion detection. *Computers & Security*, 52, 1–16.
- 51. Viegas, E. K., & Santin, A. O. (2019). Behavior-based feature selection for network intrusion detection. *Computers & Security*, 81, 1–11.
- 52. Wang, W., Zhu, M., Zeng, X., Ye, X., & Sheng, Y. (2017). Malware traffic classification using CNN for representation learning. *IEEE International Conference on Information Forensics and Security*, 798–803.
- 53. Xia, Y., et al. (2020). A hybrid unsupervised-supervised learning approach for network intrusion detection. *Computers & Security*, 91, 101747.
- 54. Yang, Q., Liu, Y., Chen, T., & Tong, Y. (2019). Federated machine learning: Concept and applications. ACM Transactions on Intelligent Systems and Technology, 10(2), 1-19.
- 55. Yang, Y., Zheng, K., & Wu, W. (2021). Machine learning for intelligent network security management. *Future Internet*, 13(2), 45.
- 56. Yin, C., Zhu, Y., Fei, J., & He, X. (2017). A deep learning approach for intrusion detection using recurrent neural networks. *IEEE Access*, 5, 21954–21961.
- 57. Zhang, J., Zulkernine, M., & Haque, A. (2008). Random-forests-based network intrusion detection systems. *IEEE Transactions on Systems, Man, and Cybernetics*, 38(5), 649–659.