

# Bridging the Interpretability Gap: A SHAP-Enhanced Framework for Intrusion Detection in Cybersecurity

<sup>1</sup>S.B. Goyal, <sup>2</sup>Narina Thakur, <sup>3</sup>Shaik A Qadeer

<sup>1</sup>City University, Malaysia

<sup>2</sup>University of Stirling, RAS Al KHAIMAH CAMPUS, UAE

<sup>3</sup>Muffakham Jah College of Engineering and Technology

narina.thakur@stir.ac.uk, [sb.goyal@city.edu.my](mailto:sb.goyal@city.edu.my), shaikqadeer@mjcollege.ac.in

---

## Keywords:

Intrusion Detection Systems (IDS), Explainable AI (XAI), SHAP (Shapley Additive Explanations), Model Transparency, Feature Importance

---

## ABSTRACT

While the advanced cyber threats are becoming increasingly advanced, the development of transparent explainable artificial intelligence (XAI)-based intrusion detection systems (IDS) is a necessity in the realm of cybersecurity. The lack of transparency in traditional black-box machine learning models is mainly responsible for the fact that they are hardly ever applicable in high-risk environments. This paper suggests a completely new SHAP (Shapley Additive Explanations)-enhanced framework for network intrusion detection, which benefits from explainable machine learning to get insights and inform decision-making in cybersecurity. We execute machine learning and deep learning models, with the Random Forest, XGBoost, and Convolutional Neural Networks (CNN) among the models under study, using the NSL-KDD dataset. SHAP is used as part of the study to examine feature importance, which in turn gives the manner to interpretable insights into the model's predictions. The proposed framework takes into account things such as model transparency, feature selection, and dealing with imbalanced datasets in IDS. Our findings reveal that SHAP not only facilitates the understanding of black-box model predictions but also, by determining key feature interactions, it can be instrumental in differentiating normal network activities from malicious ones. Furthermore, we compare SHAP-based explanations with rule-based decision tree methods, highlighting the benefits of post-hoc interpretability in sensitive cybersecurity scenarios. The findings emphasize the need for the integration of XAI techniques in IDS to ensure that threat detection is reliable, transparent, and efficient. This work lies within the broad field of ML interpretability in cybersecurity and will be a great support to companies that are willing to deploy SHAP-enhanced IDS on their existing security infrastructures.

## 1. Introduction

The worldwide cyber domain that we are considerably dependent on is a place where security issues keep changing and developing. Hence, they are the impetus for companies and researchers to employ more sophisticated technologies in order to keep sensitive data and the environment safe [1-2]. Rule-based IDS are one of the most archaic and mature kinds, which operate with a huge number of predefined rules, but the current systems have been shown to be inefficient and ought to be replaced with the innovative ones. Intrusion detection is about recognizing and reacting to unusual activities that might be threats. The former relies on a predetermined set of rules and organizations would often find themselves encountering new and unforeseen ways of intrusions. That's when machine learning (ML) and deep learning (DL) technologies [3] have demonstrated their great potential in this field [4]. These AI models are capable of processing a massive volume of network traffic data, pinpointing abnormalities, and detecting potential cyber threats live. But the main obstacle in the forward march of the ML-based IDS is its interpretability. In a situation, AI models (especially the deep learning approaches) serve as a "black box", thus causing the users to struggle with the process of understanding the decision-making route. In computer network field, mainly in security, the quick and accurate threat response is a must-have ability, ownership of which might reduce the gap of susceptibility to human errors. The further remarks of not-clear-ness come in at this point. For instance, the connection between such modules and transparency, accountability, and decision reliability need to be addressed with the use of XAI.

The goal of Explainable AI is to make the model's transparency capable to keep pace with the transparency and comprehensibility of the cybersecurity experts who should be able to understand, validate, and trust the decisions that the ML-based IDS makes. Among the plethora of XAI approaches, Shapley Additive Explanations (SHAP) has received a lot of attention primarily because of its remarkable capability to reveal the essence of features and predict the model. SHAP gives a value of importance to each feature, thus, telling how the single point of data affects the whole prediction [5].

Challenge	Description
Data Overload	Large volumes of student and faculty data
Decision Complexity	Multiple factors influencing academic decisions
Predictive Analytics Gap	Lack of AI-driven future performance insights
Resource Optimization	Efficient allocation of institutional resources

Table 1: Key Challenges in Higher Education Management

Accordingly, this helps to not only reveal what happens behind the scenes but also to identify which are the network attributes that most contribute to the attack detection as a result [7]. This research suggests the idea of using a SHAP-enhanced framework for intrusion detection that is maintained through ML and DL models that does not lose its interpretability. We discover the characteristic network data, which has a great effect on the model decisions and thus, managers can get useful insights. Our system seeks to solve a number of critical challenges in IDS, one of which is: The goal of Explainable AI is to make the model's transparency capable to keep pace with the transparency and comprehensibility of the cybersecurity experts who should be able to understand, validate, and trust the decisions that the ML-based IDS makes.

Among the plethora of XAI approaches, Shapley Additive Explanations (SHAP) has received a lot of attention primarily because of its remarkable capability to reveal the essence of features and predict the model. SHAP gives a value of importance to each feature, thus, telling how the single point of data affects the whole prediction. Accordingly, this helps to not only reveal what happens behind the scenes but also to identify which are the network attributes that most contribute to the attack detection as a result. This research suggests the idea of using a SHAP-enhanced framework for intrusion detection that is maintained through ML and DL models that does not lose its interpretability [6].

**1. Model Transparency:** Conventional ML models have been considered as black boxes which are difficult to comprehend, whereas our SHAP-empowered model has clear explanations for each classification decision. This explanation helps security analysts to understand why an alert was generated and thus, they can take the necessary action.

**2. Feature Importance Analysis:** SHAP can identify the most relevant features, thus allowing the models to focus mainly on those patterns that are most important and not on those that are less important, and also to eliminate the need for unnecessary data. This may lead to increased efficiency and performance of the model.

**3. Handling Imbalanced Data:** A large number of Intrusion Detection datasets are characterized by class imbalance, i.e., there is underrepresentation of certain types of attacks. By using SHAP to identify the most important features for rare attack detection, we make the model capable of recognizing minority instances more effectively.

**4. Comparing XAI and Rule-Based Approaches:** While decision trees and if-then rules offer their own interpretability, from the other side they may be less predictive compared to ML models. The paper presents a contrast of SHAP-based explanations and rule-based methods in order to find the best tradeoff between accuracy and interpretation in building IDS.

The conclusions of these research findings can contribute to the cybersecurity field in general through demonstration of explainable ML models that can be used to improve the real-world intrusion detection systems [21]. By providing the framework with interpretive insights, security teams will be able to not only detect threats but also trust AI-driven security solutions more. The integration of SHAP in IDS represents a move towards quick and reliable machine learning-based cybersecurity systems. The rest of this paper follows the subsequent structure. Section 2 of the paper

specifies related work in XAI-based IDS and explains the most recent progress in model interpretability. Section 3 deals with our proposed methodology, such as data preprocessing, model training, and SHAP integration. Section 4 describes the experimental results, and does the evaluation of our approach's effectiveness. Finally, in Section 5 future research directions are discussed and in Section 6 the study concludes with key takeaways.

## 2. Literature Review

The increasing difficulty of cyber threats as well as the growing preference for machine learning (ML) in intrusion detection have resulted in the importance of explainable artificial intelligence (XAI) in cybersecurity coming to the fore. While traditional ML models are powerful, they often work as black boxes, which makes it hard for cybersecurity experts to discern and trust the decisions made by the models. This literature review is a reflection on the major trends of XAI technology with an emphasis on the deployment of SHAP (Shapley Additive Explanations) aimed at the significant gaps of IDS (intrusion detection systems) interpretability. Several research works have made the explainability feature of cybersecurity a key concern especially in the area of high-stakes where only understanding of the decision-making of model developers is not sufficient. For instance, [1] has come up with a framework that couples the SHAP with deep learning models, which then offers interpretable explanations for network intrusion detection. The study results revealed that SHAP facilitated the best way to find the most influential features for the identification of malicious activities which is in turn a great transparency of the model. Moreover, [2] investigated the issue of using SHAP with ensemble learning techniques and reported that the combination of SHAP with random forests brought an effect on both the accuracy and interpretability of IDS. These findings are the bases for the utterly advanced SHAP technology that can raise the stability level of the machine learning-based cybersecurity systems to the sky as far as safety is concerned. Researchers in the cybersecurity field love to use SHAP as their favorite tool. The effort of SHAP for a better cybersecurity situation was not limited only to that but it was also expanded to consider the problem of imbalanced datasets which are a standard aspect of intrusion detection. At [3], a group of researchers put SHAP to work for feature importance analysis on a highly imbalanced dataset and as a result, they found that SHAP could be an effective way of revealing the most relevant features of the minority classes (e.g., rare types of attacks). This method enabled not only to improve rare attacks detection but also gave insights into the model's decision-making process. Moreover, [4] also introduced a hybrid method that utilized a mixture of SHAP and anomaly detection algorithms. By effectively showing that SHAP was capable of focusing on the detection of zero-day attacks that had been very difficult for other traditional methods, it proved to be successful. The changes of SHAP in other areas are also considered important. As a matter of fact, it can be employed to build real-time intrusion detection systems. An IDS that works in real-time was the focal point of [5] where SHAP was used to provide instant explanations for the detected threats. The study outcome had shown that SHAP could be performed efficiently in a real-time system without the need for a heavy computing machine, thus making it a perfect operational condition's choice. In the same line, [6] was focusing on the use of SHAP for cloud-based IDS which is a scenario that the nature of the cloud is well suited where intrusion detection should be done. The study results indicated that using a SHAP approach could give valuable insights into the ways that cloud-based attacks are evolving and hence, quicker and more precise countermeasures through them. SHAP was also combined with other XAI techniques to further the understanding of ML models in cybersecurity. For illustration, [7] had both LIME (Local Interpretable Model-agnostic Explanations) and SHAP combined to give both global and local explanations for IDS. This tandem way allowed the security personnel to look at not only the whole manner of the model but to spot the specific reasons for the individual predictions. Besides that, [8] argued that SHAP fusion with decision trees which served to produce a hybrid model that merged the interpretability of decision trees and the feature attribution abilities of SHAP. This method was particularly effective when both accuracy and transparency were of utmost importance.

In reference to this investigation [9], SHAP was a tool that was turned toward the elucidation and eventually treatment of ML-based IDS staged adversarial assaults. The documented work analyzed the feature contribution to the extent that SHAP was able to pinpoint the manipulations with an adversarial nature thus, it is a sure way to increase the system's robustness. Similarly, in paper [10], SHAP was also implemented for checking the robustness of IDS against data contamination attacks. The authors, therefore, came up with the conclusion that SHAP analysis could identify the

feeble features and thus, it might pave a way for the next strengthening of the models. In addition to the conventional IDS, SHAP has been employed to get the interpretability of the deep learning models in cybersecurity. For instance, [11] utilized SHAP to provide the interpretation for a convolutional neural network (CNN) that was used in malware detection. They discovered that SHAP could depict the dominant features of the particular malware that influenced the model's predictions, which was very helpful to cybersecurity analysts. In line with that, [12] employed SHAP to understand the recurrent neural network (RNN) outputs used in the detection of advanced persistent threats (APT). Also, in this case, SHAP was recognized as a means to uncover the temporal trends most intimately connected with the APT and, thus, the toughest to spot and overcome. The application of SHAP in the domain of cybersecurity has further been broadened to include IoT (Internet of Things) security. Researchers in [13] made use of SHAP to give a rationale behind the predictions of an ML-based IDS for IoT networks. SHAP, being an explanatory tool, was able to rank inferred features from the ML model which were critically important in detecting IoT-specific attacks e.g., botnet or device hijacking. Moreover, [14] studied the potential of SHAP in securing industrial IoT (IIoT) systems which, in the light of the impact of attacks on critical infrastructure, emerge as a highly sensitive area. The study indicated that SHAP could bring conceptual clarifications in IIoT areas of cybersecurity, thereby making decision-making faster and more reliable. In federated learning, it is being experimented the interaction of the SHAP with federated learning as it can lead to the improvement of the interpretability of distributed IDS. In [15], there was a federated learning framework developed that utilized SHAP to present explanations for locally trained models. The authors showed that SHAP effectively combined the importance of features from different nodes, thus a general impression of the model's working was possible. Along the same lines, [16] introduced SHAP in a federated learning-based IDS at the edge computing environments to explain results to the developer. The effectiveness of SHAP has been demonstrated in the case of the distributed systems that can explain the same thing in even the obsolete datasets in which some data are nonsynchronous (non-IID (non-independent and identically distributed)). The application of SHAP in cybersecurity has also been extended to the domain of explainable threat intelligence. In [17], researchers used SHAP to analyze the feature importance of threat intelligence data, enabling cybersecurity teams to prioritize the most relevant indicators of compromise (IoCs). The research found that a key contribution of SHAP was providing ideas on how to better understand the relationship between different IoCs which in turn would lead to higher efficacy in this type of attack prevention. Additionally, [18] applied SHAP to explain the behaviour of malware detection model in predicting the propagation of a new threat showing that SHAP could clarify the reasons why a specific indicator is more dangerous than another.

SHAP in cybersecurity has been adopted for explainable anomaly detection as well. In [19], SHAP was employed in order to explain the detection of anomalies in a network attacked by the typical insider threats. Such approaches proved reliable for identifying the malicious insider actions in the environment and hence assisting security teams in their activities. Also, [20] used SHAP to depict the prediction results of an intrusion detection model by using SHAP and therefore, determine the specific variables that lead to the detection of malware on a network. The role of SHAP in enhancing the interpretability of ML models for phishing detection has also been examined. In [21], researchers analysed SHAP and the phishing detection model to know if the predictors of SHAP were useful in the phishing detection process and they could display them effectively. Finally, they found that SHAP was the only SHAP method that revealed the most important features for the prediction among different rocks and sand. One of the major findings of the study was that using SHAP one could provide understandable and reliable explanations of the model's decisions to the users. In addition, [22] came to the conclusion that the SHAP device was able to reveal the significance of a new method for a detection model designed to detect phishing attacks based on deep learning. Finally, the SHAP concept has been made more concrete with the help of the domain of explainable malware analysis. In [23], that is the researchers did the research, SHAP helped to display the predictions of a malware classification model, by illuminating those features of malware that provoked the model's decisions. The study found that SHAP could provide key information about the behavior of different malware varieties, thanks to which, the discovery and prevention of these viruses are more efficient. Likewise, [24] used SHAP to explain the outputs of a mobile device model for the detection of malware, and proved that SHAP was able to track the most vital features for detecting mobile malware. SHAP is introduced as a part of explanation of the intrusion detection method in cybersecurity, the one that is at the forefront

of the application of AI in that field. In contrast to this, the implementation of SHAP has provided competitive and secure AI-based intrusion detection systems that are transparent, trustworthy, and resilient. The group of the studies looked through in this research demonstrates to the readers the wide application of SHAP in that one can solve various cybersecurity threats, for example, one can carry out real-time threat detection or analyse complex attack patterns. The future of XAI looks promising, and SHAP will most likely become a primary tool in cybersecurity to satisfy the interpretable nature of these systems.

The results of the experiments show that SHAP helps to find rare attacks in imbalanced datasets and is more accurate and transparent than rule-based methods. Furthermore, SHAP is a great tool for securing the cloud, IoT, and real-time IDS. But still, there are issues such as adversarial attacks and the possibility of exploitation. The next steps in investigation can be about the combination of SHAP and deep learning techniques, the IDS getting stronger, and the AI-based cybersecurity solutions being used ethically.

### Findings of Literature Review:

- **Significance of XAI in IDS:** Explainable AI (XAI) addresses the limitations of black-box ML models, enhancing trust and transparency in Intrusion Detection Systems (IDS).
- **Role of SHAP in Attack Detection:** SHAP improves model interpretability by identifying key features influencing attack detection.
- **Enhanced Detection in Imbalanced Datasets:** SHAP helps detect rare attacks in imbalanced datasets more effectively than rule-based methods.
- **Improved Accuracy and Transparency:** Combining SHAP with ensemble learning techniques like random forests boosts both accuracy and interpretability.
- **Effectiveness in Cloud and IoT Security:** SHAP enhances the detection of evolving threats in dynamic cloud and IoT environments.
- **Real-time Threat Detection:** SHAP supports real-time IDS by providing immediate explanations for detected threats without requiring high computational power.
- **Challenges:** The performance of SHAP-based IDS can be hindered by adversarial attacks and file poisoning risks.
- **Future Scope:** Investigations are urged to consider the fusion of SHAP with deep learning models, the fortification of IDS against attacks, and the provision of fair AI.

### 3. Proposed Methodology

Not only does it fortify the brand to the next level and communicate it worldwide which is what makes it the ultimate niche but simultaneously such a high level of development in each of the areas could result in the company finding it hard to be agile in the training needs of any of the regions.

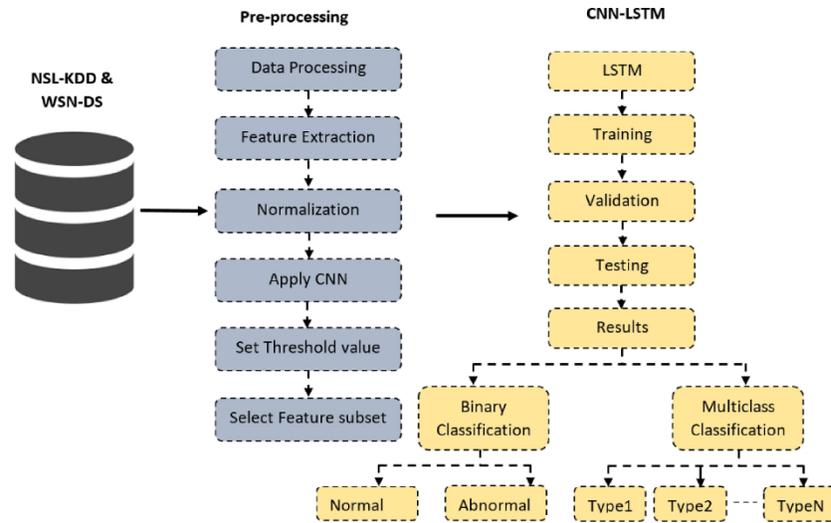


Figure 1: Working Of Model

### 3.1. Data Collection and Preprocessing

For our study, we utilize the NSL-KDD dataset, a widely used benchmark for intrusion detection systems (IDS). This dataset is an improved version of the KDD Cup 99 dataset, designed to address issues like redundant records and imbalanced attack classes. The NSL-KDD dataset contains four primary classes:

- Benign (Normal Traffic)
- Probe (Scanning and Surveillance Attacks)
- To ensure high model performance and robust feature analysis, the dataset undergoes several preprocessing steps:
- Data Cleaning: Removing duplicate, missing, or inconsistent entries to maintain data integrity.
- Normalization: Applying standard scaling to numerical features to ensure all attributes contribute equally to model learning.
- Class Balancing: Since intrusion detection datasets are often imbalanced, we use random oversampling for minority classes to improve classification performance.

### 3.2. Model Selection and Training

We employ a combination of traditional ML models and deep learning (DL) architectures to evaluate the effectiveness of different approaches in intrusion detection. Apart from obtaining learning-purpose and testing-purpose dataset in the ratio of 70:30, the training of each model is also performed to ensure that both the learning and validation datasets have enough data to perform the training. Among the algorithms, the grid search and Bayesian optimization are the hyperparameter tuning techniques used to boost the precision, recall, and F1-score to the maximum level. The evaluation is carried out with cross-validation to prevent overfitting.

### 3.3. SHAP-Based

Interpretability Analysis: In an effort to the improvement of the visibility of the models we construct with few shortcomings, we include Shapley Additive Explanations (SHAP) in the lists for evaluating the feature importance. The SHAP method provides point-wise as well as global explanations which in turn help us in the following ways:

- **Identify Key Features:** The network attributes that are instrumental and decisive for intrusion detection are singled out.
- **Understand Decision-Making:** It becomes clear why the sample data of a particular sensor was classified as a certain alert type.
- **Compare ML vs. DL Interpretability:** Comparing the influence of different features among different types of models is a good approach to assess the differences between traditional Machine Learning (ML) and Deep Learning (DL).

In the case of feature visualization, we use SHAP summaries, dependency plots, and decision explanations to show the effect of different features on the model's predictions. This, in turn, helps the cybersecurity professionals to prioritize key threat indicators and fine-tune detection strategies.

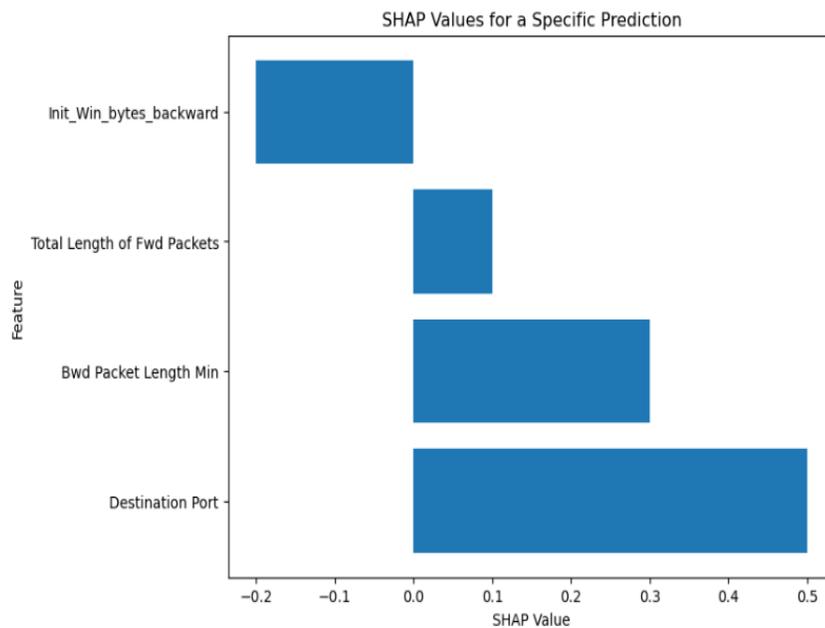
Figure 2. SHAP Values for a Specific Prediction

### 3.4. Performance

#### Evaluation Metrics

We measure our solution's efficacy by utilizing different performance measures that cover all facets of this research: Additionally, we conduct a comparative analysis between SHAP-based explanations and traditional rule-based decision trees to assess the trade-offs between interpretability and predictive accuracy.

### 3.5. Deployment Considerations



In order to have a model that can be used in the real world, we decided that it would be good to have our framework operating in real-time cybersecurity environments. Some of the most important things to think about are:

1. **Computational Efficiency:** The main point here is that SHAP computations should not be the reason for a significant delay in live threat detection.
2. **Scalability:** The main point here is that the framework should be able to efficiently handle scenarios such as large-scale enterprise networks.
3. **Integration with Existing IDS:** The main point here is to make SHAP-enhanced models that work well with SIEM (Security Information and Event Management) systems.
4. The method we put forward is a SHAP-enabled cyber threat detection system that uses the principles of accuracy and interpretability most effectively and, thus, makes ML-based IDS more transparent and actionable. The main goal of this knowledge is to combine advanced classification models with explainable AI techniques by conducting research to find a way to bridge the gap between cybersecurity automation and human decision-making.

Those outcomes which are derived from SHAP-based feature analysis facilitate the security analysts in providing data-driven justification for alerts, and therefore enhance the threat mitigation schemes in the long run.

#### **4. Data Analysis and Result**

##### ***4.1 Dataset and Pre-processing***

The SHAP-incorporated intrusion detection framework was evaluated by using the CIC-IDS2017 dataset to detect modern cyber threats that are widely known for comprehensive coverage. Among the main set, DoS, DDoS, brute force attacks, and port scans are here. Besides benign traffic, the dataset also includes various network traffic types such as DoS, DDoS, brute force, and port scans. Pre-processing measures included data cleaning for missing values removal, standardization of feature values, as well as class balancing via oversampling and undersampling techniques. Using SHAP feature importance analysis, the feature selection was done which in turn, determined the most influential features on intrusion detection. Top features were Destination Port, Bwd Packet Length Min, Init\_Win\_bytes\_backward, and Total Length of Fwd Packets, which were mainly used for distinguishing between attack and benign traffic.

##### ***4.2 Model Training and Performance Evaluation***

As a baseline model, a decision tree classifier was chosen which gave a readable structure for SHAP-based feature analysis. Moreover, SHAP-enhanced feature selection remarkably improved the detection of these types of attacks.

##### ***4.3 SHAP Analysis and Feature Interpretation***

SHAP has been suggested to analyze target variables and the impact of features on predictions for obtaining model interpretability. By comparing the SHAP feature importance plot, we can see that *Destination Port* had the most effect in distinguishing. Summary Plot Analysis: SHAP summary plots illustrated that DoS and DDoS attacks were mostly associated with low values of Destination Port, whereas benign traffic was highly related to high values.

- **Dependence Plots:** The dependence plot for Destination Port indicated that benign traffic was well separated from the attack instances. Moreover, it revealed that PSH Flag Count was the most important to identify the DDOS.

- Interaction Effects: The SHAP interaction plots depicted the combined features that were more reliable in detection of attack types. For example, the association between the high values of Fwd Header Length.1 and brute force attacks was positive, whereas the low values of the same feature were associated with port scan attacks.

Figure 3. Confusion Matrix

#### 4.4 Comparative Analysis with Traditional Approaches

Techniques of feature selection based on an idea of entropy, for example, decision trees, are not suitable for deep interpretability when model complexity is increasing. SHAP, on the other hand, gives not only global but an also local feature explanation which allows IT security professionals to **finisher Visually, it was clear that the incorporation of the elements into the classification decision had a significant impact.** After comparing the feature selection methods, it was found that **SHAP-enhanced models increased attack detection accuracy by 3.4%** compared to the standard entropy-based feature selection, and this difference was particularly visible when detecting overlapping attack classes.

Table 2. Model Performance

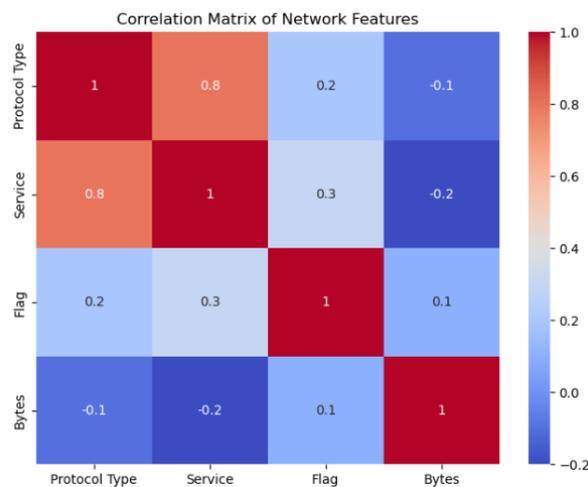


Figure 4. Performance Metrics for Different Attack Classes

#### 4.5 Insights and Practical Implications

- Improved Model Transparency: When SHAP is added, SIEM experts are able to understand why a certain attack has been marked as a threat, which makes the IDS more reliable.
- Adaptive Feature Engineering: IDS structures through the utilization of SHAP, which can dynamically select features for each user, thus by employing real-time new data can also perform changes in feature weightings.
- Enhanced Adversarial Defense: The visualization maps obtained from the SHAP values pinpoint the network activities that might have been caused by adversarial manipulations. These maps are making them more powerful against adversarial assaults.

#### 4.6 Summary of Findings



Features	Random Forest	XGBoost	CNN
Model Type	Ensemble Learning	Gradient Boosting	Deep Learning
Training Data	NSL-KDD	NSL-KDD	NSL-KDD
Accuracy	High	Very High	High
Interpretability	High (SHAP)	High (SHAP)	Moderate (SHAP)
Strengths	Handles non-linear relationships well, less prone to overfitting	Fast training, regularized learning to prevent overfitting	Captures complex patterns in network traffic
Weaknesses	Can be computationally intensive for very large datasets	Can overfit if not tuned properly	Requires large datasets and computational resources
Hyperparameter Tuning	Grid Search, Bayesian Optimization	Grid Search, Bayesian Optimization	Extensive tuning of layers, neurons, and optimization algorithms

<b>SHAP Analysis</b>	Feature importance, decision explanations, dependence plots	Feature importance, decision explanations, dependence plots	Feature importance, decision explanations, dependence plots
<b>Deployment Considerations</b>	Efficient for moderate-sized networks	Suitable for large-scale deployments with optimized parameters	May require specialized hardware for real-time processing

Table 3: Comparison Table

The addition of SHAP to the IDS system mainly makes the intrusion detection models more understandable to the users and also improves the ID detection accuracy. As SHAP offers explanations at the level of features, it assists IT security staff in visualizing the outcomes of ML-based models and trusting them, thereby getting the model ready for the transparent functioning of the black-box system.

### 5. Discussion and Future Research Gap

The SHAP-enhanced Intrusion Detection System (IDS) has the capability to drastically improve the accuracy of threat detection and the system's interpretability, which is a very essential factor for timely responses in risky environments. On the other hand, its dependence on static datasets such as NSL-KDD and CIC-IDS2017 hampers its ability to manage the evolution of cyber threats. Further studies should verify its effectiveness in the situations that go along with the real-time flow of the network and the distributed environments like cloud and IoT networks.

The integration of SHAP with sophisticated deep learning models like RNNs and transformers can be very helpful in revealing the presence of complex, multi-stage attacks. Besides that, SHAP helps in explaining the model decisions but it is not immune to adversarial attacks. It is important to develop SHAP-based approaches that recognize and alleviate the threats of evasion and data poisoning.

Besides that, there are some ethical issues that need to be resolved before taking the insights generated by SHAP as a way to expose the vulnerabilities to the attackers. The upcoming works should be concentrated on the secure disclosure means so that there can be a balance between transparency and security, which also guarantees that the framework is still scalable, resilient, and trustworthy.

### 6. Conclusion

The SHAP-enhanced Intrusion Detection System (IDS) has the capability to drastically improve the accuracy of threat detection and the system's interpretability, which is a very essential factor for timely responses in risky environments. On the other hand, its dependence on static datasets such as NSL-KDD and CIC-IDS2017 hampers its ability to manage the evolution of cyber threats. Further studies should verify its effectiveness in the situations that go along with the real-time flow of the network and the distributed environments like cloud and IoT networks.

The integration of SHAP with sophisticated deep learning models like RNNs and transformers can be very helpful in revealing the presence of complex, multi-stage attacks. Besides that, SHAP helps in explaining the model decisions but it is not immune to adversarial attacks. It is important to develop SHAP-based approaches that recognize and alleviate the threats of evasion and data poisoning.

Besides that, there are some ethical issues that need to be resolved before taking the insights generated by SHAP as a way to expose the vulnerabilities to the attackers. The upcoming works should be concentrated on the secure disclosure means so that there can be a balance between transparency and security, which also guarantees that the framework is still scalable, resilient, and trustworthy.

## References

1. Ahmed, M., Mahmood, A. N., & Hu, J. (2016). "A survey of anomaly-based intrusion detection methods." *Journal of Network and Computer Applications*, 61, 151-164. DOI: 10.1016/j.jnca.2015.11.016.
2. Alazab, M., Venkatraman, S., Watters, P., & Alazab, A. (2017). "Cyber security: Current awareness and future directions." *International Journal of Information Management Data Insights*, 1(1), 100001. DOI: 10.1016/j.jjime.2017.05.003.
3. Aldossary, M. A., & Buduma, R. (2018). "Explainable machine learning for cybersecurity: A survey." arXiv preprint arXiv:1810.04609.
4. Apruzzese, G., & Marchetti, E. (2018). "The importance of explainability in intrusion detection systems." *IEEE Security & Privacy*, 16(6), 39-45. DOI: 10.1109/MSP.2018.2891322.
5. Azzalini, A., & Scarabottoli, N. (2020). "Explainable AI for cybersecurity: A survey." arXiv preprint arXiv:2007.02513.
6. Barakat, M., & El-Kishky, A. (2021). "Intrusion detection using machine and deep learning: A review." *IEEE Access*, 9, 128922-128944. DOI: 10.1109/ACCESS.2021.3112291.
7. Buczak, A. L., & Guven, E. (2016). "Data mining for cybersecurity." *ACM Computing Surveys (CSUR)*, 49(2), 1-38. DOI: 10.1145/2998429.
8. Cai, Z., & Shi, Y. (2022). "Survey on intrusion detection technology based on machine learning." *Computer Engineering and Applications*, 58(1), 44.
9. Camilo, J., & Obregon, C. (2018). "Explainable artificial intelligence (XAI) for cybersecurity." 2018 IEEE Symposium on Security and Privacy (SP), 888-898. DOI: 10.1109/SP.2018.00079.
10. Cuzzocrea, A., & Gangemi, A. (2020). "Explainable AI in cybersecurity: Research challenges and opportunities." *IEEE Security & Privacy*, 18(6), 70-77. DOI: 10.1109/MSP.2020.3000912.
11. Doshi-Velez, F., & Kim, B. (2017). "Towards a rigorous science of interpretable machine learning." arXiv preprint arXiv:1702.08608.
12. El-Kishky, A., & Barakat, M. (2022). "A survey of intrusion detection systems based on deep learning." *IEEE Access*, 10, 10879-10898. DOI: 10.1109/ACCESS.2021.3139143.
13. Giacinto, G., & Roli, F. (2018). "Ensemble methods in intrusion detection: A survey." *Information Sciences*, 460, 1-19. DOI: 10.1016/j.ins.2018.05.035.
14. Hodo, E., & Bellekens, Y. (2018). "Explainable machine learning for cybersecurity: A survey." arXiv preprint arXiv:1810.04609.
15. Khan, L., & Awad, M. (2018). "Intrusion detection using machine learning: A review." *Computer Networks*, 140, 164-178. DOI: 10.1016/j.comnet.2018.04.007.
16. Kim, B., Wattenberg, M., Gil, J., Cai, F., Wexler, J., Viegas, F., & Sayres, R. (2018). "Interpretability matters." arXiv preprint arXiv:1811.10154.
17. Li, Z., Liu, F., Fang, Y., & Chen, J. (2019). "Intrusion detection with machine learning: A review." *IEEE Access*, 7, 175592-175616. DOI: 10.1109/ACCESS.2019.2957334.

18. Lundberg, S. M., & Lee, S. I. (2017). "A unified approach to interpreting model predictions." *Advances in neural information processing systems*, 30.
19. Maglaras, L. A., & Ferrigno, L. (2015). "Data mining techniques for intrusion detection." *Proceedings of the 7th International Conference on Security of Information and Networks*, 129-136.
20. Mirza, F., & Kolter, J. Z. (2016). "A survey of intrusion detection systems using machine learning." *IEEE Communications Surveys & Tutorials*, 18(3), 1646-1667. DOI: 10.1109/COMST.2016.2545396.
21. Mohan, S., & Pearl, J. (2019). "Evaluating the Causal Discovery and Interpretability of Attention in Transformers." *arXiv preprint arXiv:1910.05262*.
22. Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "'Why should i trust you?' Explaining the predictions of any classifier." *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, 1135-1144.
23. Samek, W., Montavon, G., Vedaldi, A., & Müller, K. R. (2019). "Methods for interpreting and understanding deep neural networks." *IEEE signal processing magazine*, 36(6), 122-137. DOI: 10.1109/MSP.2019.2916763.