

Hypertuning Random Forest for Enhancing Cyber Security in Industrial Power Electronics

MD. Shameem Ahammed

Electrical and Electronic Engineering
Jatiya Kabi Kazi Nazrul Islam
University

Trishal, Mymensingh 2224, Bangladesh
shameem_19102913@jkkniu.edu.bd

Md. Mahbubur Rahman

Electrical and Electronic Engineering
Jatiya Kabi Kazi Nazrul Islam
University

Trishal, Mymensingh 2224, Bangladesh
mahbubur@jkkniu.edu.bd

Sourav Sarker

Electrical and Electronic Engineering
Jatiya Kabi Kazi Nazrul Islam University
Trishal, Mymensingh 2224, Bangladesh
souravsarker50.eee@gmail.com

Mobussir Nur Sium

Electrical and Electronic Engineering
Jatiya Kabi Kazi Nazrul Islam
University

Trishal, Mymensingh 2224, Bangladesh
mobussir_19102924@jkkniu.edu.bd

Faria Ferdous

Electrical and Electronic Engineering
Jatiya Kabi Kazi Nazrul Islam University

Trishal, Mymensingh 2224, Bangladesh
f.ferdouscw@gmail.com

Fatema Tanbin Mim

Electrical and Electronic Engineering
Jatiya Kabi Kazi Nazrul Islam University
Trishal, Mymensingh 2224, Bangladesh
tanbinmim2@gmail.com

Md. Shake Farid Uddin

Assistant Manager(Software Engineer)
Nuclear Power Plant Company
Bangladesh Limited.
Dhaka, Bangladesh
sfuddin.iit@gmail.com

Arit Sarker

Assistant Superintendent of police
Bangladesh Police.
Dhaka, Bangladesh
sarkeraritarup@gmail.com

Abstract—Industrial Power Electronics (IPE) systems are increasingly integrated with communication technology, creating new vulnerabilities for cyber-attacks. In most cases, traditional security measures are insufficient to deal with the complexity and evolution of attacks that target critical infrastructure. This paper presents a machine learning-based framework for detecting cyberattacks in IPE systems, leveraging the Random Forest algorithm to classify anomalous network traffic. Using the NSL-KDD dataset, which encompasses multiple attack types such as *Denial of Service* (DoS), *Probe*, *Remote to Local* (R2L), and *User to Root* (U2R) attacks, we achieved a detection accuracy of 99.61%. Our approach highlights the potential of Random Forest models in effectively identifying intrusion attempts and enhancing the resilience of industrial systems against cyber threats. The high accuracy underscores the effectiveness of feature-rich datasets and the importance of machine learning in modern industrial cybersecurity solutions.

Keywords— Industrial Power Electronics, cyber-attacks, ML, Random Forest.

I. INTRODUCTION

Industrial Power Electronics (IPE) systems are at the core of energy management and industrial automation, facilitating efficient power conversion, control, and distribution across diverse sectors. IPE systems are becoming increasingly integrated with advanced communication technologies and *Industrial Internet of Things* (IIoT) networks, which expose them to a growing number of cybersecurity threats[1][2]. Often used in power grids and manufacturing plants, these systems are vulnerable to sophisticated cyberattacks because of their critical nature. The impact of such attacks can be severe, resulting in operational disruption, financial losses, and, in some cases, safety concerns[3].

In recent years, the threat landscape has evolved, with cyberattacks becoming more frequent and complex. Attack vectors such as *Denial of Service* (DoS), *Probe*, *Remote to Local* (R2L), and *User to Root* (U2R) attacks are among the most prevalent methods used to compromise industrial networks. DoS attacks, for example, overwhelm system resources, rendering services unavailable, while R2L and U2R attacks exploit system vulnerabilities to gain unauthorized access and escalate privileges[4]. It is important to recognize that these types of attacks can severely compromise the operational integrity of industrial systems if they are not detected and mitigated quickly.

The traditional *intrusion detection system* (IDS), which relies heavily on signature-based techniques, has proven inadequate in handling the dynamic and sophisticated nature of modern cyber threats. There is a limitation to the ability of these systems to detect zero-day attacks or emerging threats that do not conform to conventional attack patterns. Consequently, there is a pressing need for more adaptive and intelligent cybersecurity solutions capable of identifying both known and unknown threats in real-time[5].

Machine learning (ML) has emerged as a powerful tool for enhancing cybersecurity in industrial environments. ML-based IDS can learn from historical data and identify anomalous behaviour that may indicate a cyberattack[6][7]. Among the various machine learning algorithms, the Random Forest classifier stands out due to its robustness, scalability, and ability to handle high-dimensional data. By training the model on a comprehensive dataset that captures different types of network behaviours, it becomes possible to detect malicious activities with high accuracy.

In this paper, we propose a machine learning framework based on the *Hypertuned Random Forest* (HRF) model to detect anomalies in industrial power electronics systems. We

use the NSL-KDD dataset, a well-established benchmark for evaluating intrusion detection systems. This dataset contains 41 features that describe various aspects of network traffic, including basic, content-based, and time-based attributes, and encompasses multiple attack types. The HRF model was trained and tested using this dataset, achieving a high accuracy of 99.61% in classifying normal and anomalous conditions. This research demonstrates that with proper feature selection and model optimization, Random Forest classifiers can significantly enhance the cybersecurity posture of industrial systems.

Industrial systems can be better prepared to identify and react to cyber-attacks in real-time, protecting the security and dependability of vital infrastructure, by utilizing the insights acquired from this effort. The outcomes of this study show the possibility for incorporating these models into real-world settings and open the door for additional research into machine learning approaches in industrial cybersecurity.

The rest of the paper are organized as follows: Section II provides a detailed explanation of the proposed model, including dataset, data pre-processing and feature engineering, Random Forest algorithm, Hyperparameter tuning using RandomizedSearchCV, model training and evaluation. Section III discusses the results including Confusion matrix and ROC curve. Finally, Section IV conclude the paper with some future works.

II. PROPOSED MODEL

The proposed approach applies a Random Forest classifier to detect cyber anomalies in IPE systems. We evaluated the model using the NSL-KDD dataset, which contains labelled network traffic data with a mix of normal and attack instances. This section outlines the dataset, preprocessing steps, model selection, and hyperparameter tuning strategies used in the study.

A. Dataset

The dataset is collected from Kaggle website[8]. The NSL-KDD dataset consists of 41 features that describe various network behaviours. These features fall into three main categories.

Basic Features: Characteristics directly extracted from TCP/IP connections, such as duration, protocol type, and service.

Content Features: Features based on the payload of network packets, such as num_failed_logins and is_guest_login.

Traffic-Based Features: Metrics related to traffic within a specific time window, such as error_rate (percentage of SYN errors) and srv_count (number of connections to the same service).

The dataset includes attacks categorized as *Denial of Service* (DoS), *Probe*, *Remote to Local* (R2L), and *User to*

	Duration	Protocol Type	Service	Flag	Srv Bytes	Dst Bytes	Land	Wrong Fragment	Urgent	Hot	...	Dst Host Srv Count	Dst Host Same Srv Rate	Dst Host Diff Srv Rate	Dst Host Srv Error Rate	Dst Host Same Srv Error Rate	Dst Host Diff Srv Error Rate	Dst Host Srv Error Rate	Dst Host Same Srv Error Rate	Dst Host Diff Srv Error Rate	target
148507	0	1	24	9	334	1600	0	0	0	0	...	255	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
148508	0	1	49	5	0	0	0	0	0	0	...	12	0.05	0.06	0.00	0.00	1.00	1.0	0.00	0.0	1
148509	0	1	54	9	2233	365	0	0	0	0	...	2	1.00	0.00	1.00	1.00	0.00	0.0	0.00	0.0	0
148510	0	1	49	5	0	0	0	0	0	0	...	13	0.05	0.07	0.00	0.00	1.00	1.0	0.00	0.0	1
148511	0	1	24	9	359	375	0	0	0	0	...	255	1.00	0.00	0.33	0.04	0.33	0.00	0.00	0.0	0
148512	0	1	49	5	0	0	0	0	0	0	...	25	0.10	0.06	0.00	0.00	1.00	1.0	0.00	0.0	1
148513	8	2	49	9	105	145	0	0	0	0	...	244	0.96	0.01	0.01	0.00	0.00	0.0	0.00	0.0	0
148514	0	1	54	9	2231	384	0	0	0	0	...	30	0.12	0.06	0.00	0.00	0.72	0.0	0.01	0.0	0

Fig 2: Dataset

Root (U2R), along with normal traffic instances. Fig 1 shows tail part of the dataset.

It is possible to assess the model's capacity to distinguish between malicious and authentic traffic according to this

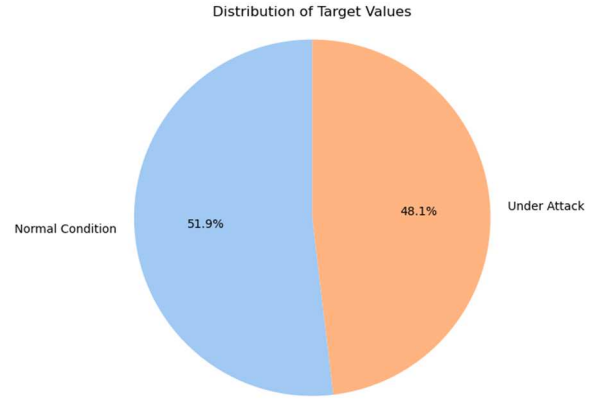


Fig 1: Distribution of target values

variety of attack types. The target values' distribution within the dataset is displayed in Fig 2

B. Preprocessing and Feature Engineering

The preprocessing phase was crucial to ensure the integrity and uniformity of the dataset:

Handling Missing Values: Any missing data in the dataset was handled by imputation with the median value or removal of problematic instances.

Encoding Categorical Features: Categorical features like protocol type and service were one-hot encoded to convert them into numerical formats suitable for the machine learning model using label encoding process.

Standardization: Numerical features were scaled using StandardScaler to bring all values into a uniform range, improving the model's performance during training.

The Correlation heatmap of relevant features and target values shows in the Fig 3.

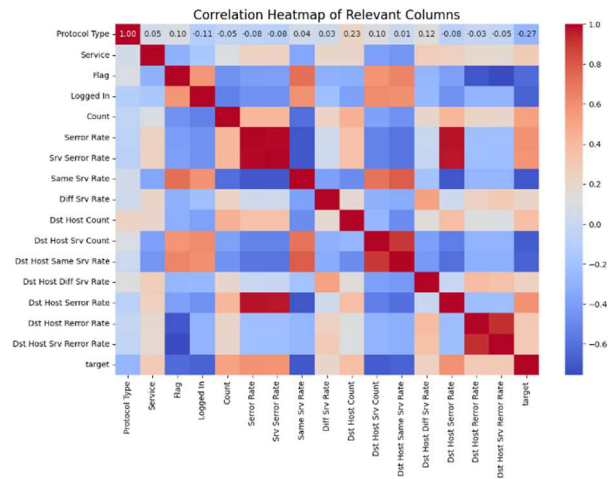


Fig 3: Correlation heatmap

It shows that some features have positive correlation, negative correlation and weak or no correlation with the target values. When a change in a feature is linked to an increase in the target value, it has a positive correlation with the target; when a change in the feature is linked to a drop in

the target value, it has a negative correlation with the target; and when a change in the feature has little to no effect on the target, it has weak or no correlation with the target. The feature is not informative for making predictions about the target value.

- **Error-related features** (Serror Rate, Error Rate) strongly correlate with the target.
- **Host and service access patterns** (Same Srv Rate, Dst Host Count) are crucial in differentiating between normal and attack traffic.

C. Random Forest Algorithm

A Random Forest classifier was chosen for its reliability in handling large datasets and complex relationships between features. Random Forest operates by building an ensemble of decision trees and averaging their predictions, which reduces overfitting and improves model accuracy. This method is particularly suitable for anomaly detection in complex systems like IPE, where patterns of attacks are often multifaceted. The default parameters used in the model are explained below:

- **n_estimators=100**; This indicates how many trees there are in the forest. The model performs better and requires more training time when there are more trees.
- **criterion='gini'**; The function ('gini' or 'entropy') that is used to gauge the quality of a split is the criteria.
- **max_depth=None**; The trees will continue to develop until there are either no more leaves than min_samples_split samples or pure leaves. Overfitting and deep trees may result from this.
- **min_samples_split=2**; The least amount of samples needed for an internal node to divide.
- **min_samples_leaf=1**; The lowest number of samples needed at a leaf node.
- **max_features='auto'**; The quantity of features to take into account while choosing the ideal split. When it comes to categorization, auto equals $\sqrt{n_features}$.
- **bootstrap=True**; Bootstrap sampling is used to train each tree on a sample of the data with replacement. This increases model variance but reduces bias.
- **n_jobs=None**; The number of parallel jobs. Using n_jobs=-1 will use all available CPU cores.

The NSL-KDD dataset uses Random Forest because it offers high accuracy, manages imbalances in classes, is durable to noise and overfitting, effectively handles multi-dimensional data, and provides insights into the significance of individual features. Because of this, it is a very powerful algorithm for network security intrusion detection. Random Forest classifier uses multiple decision tree to make the model effective and accurate. By default RF algorithm, the model using this large dataset is achieved an accuracy of 99.58%.

Here Fig 4 shows the upper four level of the decision tree. There are three components to it, which are the root node, the decision node, and the leaf node. Root nodes are the points from which the dataset begins to split. Those nodes that are obtained after splitting a root node are known as

decision nodes, whereas the nodes in which further splitting is not possible are known as leaf nodes. We will use the impurity of our dataset as the root node which will give us the lowest impurity or lowest Gini index of our dataset.

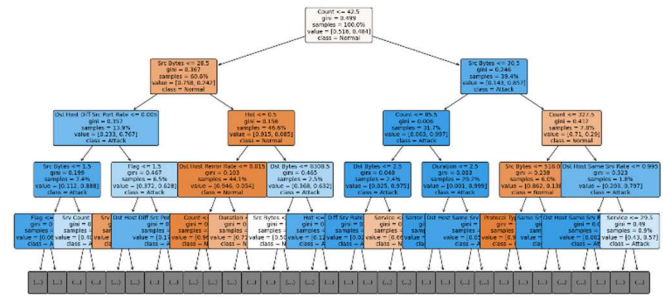


Fig 4: Decision tree of four level

Although it is higher accuracy for this detection, we will try to increase the accuracy as much as possible. So we decided to use RandomizedSearchCV to hyperparameter tuning for this purpose.

D. Hyperparameter Tuning using RandomizedSearchCV

To ensure the model reached its full potential, we used RandomizedSearchCV for hyperparameter tuning. Several important hyperparameters were tuned, including:

- **n_estimators=400**;
- **max_depth=70**;
- **min_samples_split=10**;
- **min_samples_leaf=1**;
- **max_features='auto'**;
- **bootstrap=False**;

The Randomized SearchCV method allowed for efficient exploration of the hyperparameter space, yielding the best configuration for the dataset at hand. Here some key points to using RandomizedSearchCV for this prediction.

- **Eliminates underfitting and overfitting** by optimizing important hyperparameters.
- **Finds better combinations of parameters** to prevent overfitting.
- **Improves handling of class imbalance** by adjusting weights and other variables.
- **Increases the robustness** of the model by employing cross-validation to provide more accurate performance estimation.

The model achieved an accuracy of 99.61% after using RandomizedSearchCV to adjust the hyperparameter, which is somewhat better than the default RF model.

E. Model Training and Evaluation

The dataset was split into test (30%) and training (70%). The model was tested on the test set after being trained on the training set. Performance was assessed using important metrics including F1-score, accuracy, precision, and recall. Additionally, a confusion matrix was created to offer a thorough explanation of how the model categorized different sorts of attacks.

III. RESULT AND DISCUSSION

The Hypertuned Random Forest model demonstrated excellent performance, achieving an accuracy of 99.63% on the NSL-KDD dataset. The model showed strong capability in detecting DoS and Probe attacks. However, it exhibited

slightly lower performance in identifying less frequent attack types, such as R2L and U2R attacks.

A. Confusion Matrix

The confusion matrix highlights the model's ability to correctly classify each attack type. The model performed exceptionally well in detecting DoS and Probe attacks, as these were well-represented in the dataset. However, the model struggled with R2L and U2R attacks, which had lower representation in the dataset.

Here the matrix format in Table 1.

Table 1: CONFUSION MATRIX

True Positive(TP)	False Negative(FN)
False Positive(FP)	True Negative(TN)

Where,

- TP refers to truly predict positive outcomes.
- TN refers to truly predict negative outcomes.
- FP stands for falsely predict positive outcomes
- FN stands for falsely predict negative outcomes

Here Fig 5 shows the Confusion Matrix of default Random Forest Algorithm.

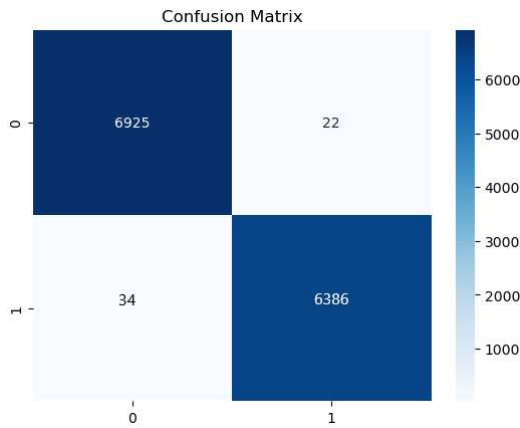


Fig 5: Confusion Matrix of Default RF

The Fig 6 shows the Confusion Matrix of our proposed model of Random Forest which is hypertuned using RandomSearchCV.

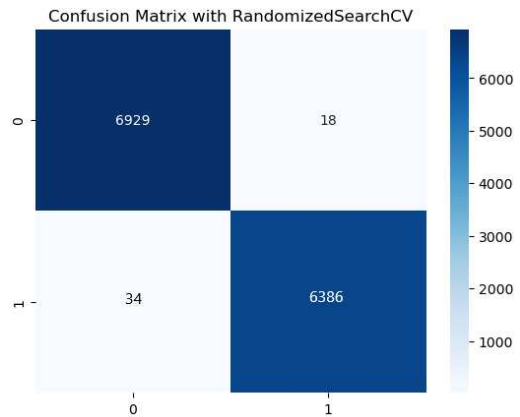


Fig 6: Confusion Matric of RF with RandomizedSearchCV

B. ROC Curve

For every attack class, the Receiver Operating Characteristic (ROC) curve was displayed to show the trade-off between true positive rates and false positive rates. The model demonstrated a high Area Under the Curve (AUC) for both DoS and Probe assaults, indicating its efficient ability to differentiate between normal and dangerous information. ROC uses two parameters:

- True Positive Rate (TPR):

$$TPR = \frac{TP}{TP+FN} \dots\dots\dots 1$$

- False Positive Rate (FPR):

$$FPR = \frac{FP}{FP+TN} \dots\dots\dots 2$$

Area Under the Curve denotes as AUC. To compute the points in an ROC curve of different classification thresholds provided by AUC. It counts the entire 2D area underneath the ROC[9]. The ROC curve (Fig 7) with AOC values for our proposed model is given below.

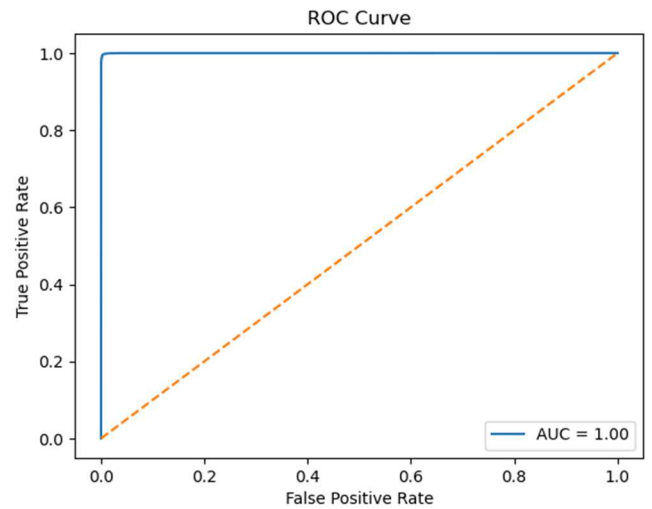


Fig 7: ROC curve

Accuracy means how close the predicted and actual values are.

$$Accuracy = \frac{TP + TN}{TP + FP + FN + TN} \dots\dots\dots 3$$

Precision measures true positive prediction among all positive predictions.

$$Precision = \frac{TP}{TP+FP} \dots\dots\dots 4$$

Recall measures the positive prediction among all actual positive items.

$$Recall = \frac{TP}{TP + FN} \dots\dots\dots 5$$

f1-score is a metric that balances precision and recall.

$$f1 - score = 2 * \frac{Precision * Recall}{Precision + Recall} \dots\dots\dots 6$$

Now, the result of default RF model and Hypertuned Random Forest means our proposed model including precision, recall, f1-score, and accuracy are given below.

Table 2: Result Analysis

Model	Precision n (%)	Recall (%)	f1- score (%)	Accuracy (%)
General RF	99.51	99.68	99.59	99.58
Hypertuned RF(Proposed Model)	99.51	99.74	99.62	99.61

The results demonstrate that Random Forest is a highly effective algorithm for anomaly detection in industrial power electronics systems. The model's ability to detect DoS and Probe attacks with high accuracy confirms its suitability for real-time deployment in critical infrastructure. However, the lower performance on R2L and U2R attacks highlights the need for further refinement, potentially through the inclusion of additional features or more advanced algorithms.

The model's high overall accuracy of 99.61% suggests that Hypertuned Random Forest can form the backbone of a robust intrusion detection system (IDS) in industrial environments. Future work should focus on improving the detection of rare attack types, perhaps by leveraging more advanced ensemble methods or deep learning techniques. Additionally, integrating this model with live IPE systems could provide real-time cybersecurity monitoring, helping to safeguard critical infrastructure from emerging cyber threats.

IV. CONCLUSION

This study presents a Hypertuned Random Forest-based machine learning framework for enhancing cybersecurity in Industrial Power Electronics systems. By training the model on the NSL-KDD dataset, we demonstrated its ability to detect network anomalies with an accuracy of 99.61%. The model is particularly effective at detecting high-frequency attack types like DoS and Probe, making it a promising solution for protecting industrial systems against cyberattacks. Ensure low-latency inference and smooth integration with *Industrial Control Systems* (ICS) by deploying the optimized machine learning model on an appropriate platform. Utilize real-time data processing methods, such as batch and asynchronous processing, to minimize delay and maximize performance. Utilize APIs and strong security mechanisms to integrate the ML model with ICS components. Implement alert and

response mechanisms to trigger appropriate actions. Future work will explore more sophisticated models and techniques to improve the detection of rarer attack types, such as R2L and U2R. Additionally, the integration of this model into real-world industrial systems offers significant potential for real-time anomaly detection and response.

REFERENCES

- [1] Gonaygunta, H., Nadella, G. S., Pramod Pawar, P., & Kumar, D. (2024). Enhancing Cybersecurity: The Development of a Flexible Deep Learning Model for Enhanced Anomaly Detection. *2024 Systems and Information Engineering Design Symposium, SIEDS 2024*, 79–84. <https://doi.org/10.1109/SIEDS61124.2024.10534661>
- [2] Khan, N., & Ammar Taqvi, S. A. (2023). Machine Learning an Intelligent Approach in Process Industries: A Perspective and Overview. *ChemBioEng Reviews*, 10(2), 195–221. <https://doi.org/10.1002/CBEN.202200030>
- [3] Mazhar, T., Irfan, H. M., Khan, S., Haq, I., Ullah, I., Iqbal, M., & Hamam, H. (2023). Analysis of Cyber Security Attacks and Its Solutions for the Smart grid Using Machine Learning and Blockchain Methods. *Future Internet*, 15(2). <https://doi.org/10.3390/fi15020083>
- [4] Mohammed, S. H., Al-Jumaily, A., Singh, M. S. J., Jimenez, V. P. G., Jaber, A. S., Hussein, Y. S., Al-Najjar, M. M. A. K., & Al-Jumeily, D. (2024). A Review on the Evaluation of Feature Selection Using Machine Learning for Cyber-Attack Detection in Smart Grid. *IEEE Access*, 12(March), 44023–44042. <https://doi.org/10.1109/ACCESS.2024.3370911>
- [5] Mughaid, A., Aljamal, M., Al-Aiash, I., Aljamal, M., Alquran, R., Alzu'bi, S., & Abutabanjeh, A. A. (2023). Enhancing Cybersecurity in SCADA IoT Systems: A Novel Machine Learning-Based Approach for Man-in-the-Middle Attack Detection. *2023 3rd Intelligent Cybersecurity Conference, ICSC 2023*, 74–79. <https://doi.org/10.1109/ICSC60084.2023.10349993>
- [6] Sahani, N., Zhu, R., Cho, J. H., & Liu, C. C. (2023). Machine Learning-based Intrusion Detection for Smart Grid Computing: A Survey. *ACM Transactions on Cyber-Physical Systems*, 7(2). <https://doi.org/10.1145/3578366>
- [7] Zhukabayeva, T., Pervez, A., Mardenov, Y., Othman, M., Karabayev, N., & Ahmad, Z. (2024). A Traffic Analysis and Node Categorization-Aware Machine Learning-Integrated Framework for Cybersecurity Intrusion Detection and Prevention of WSNs in Smart Grids. *IEEE Access*, 12(June), 91715–91733. <https://doi.org/10.1109/ACCESS.2024.3422077>
- [8] Hassan, M. (2022). *NSL-KDD Dataset for Network Anomaly Detection* [Data set]. Kaggle. <https://www.kaggle.com/datasets/hassan06/nslkdd>
- [9] Google for Developers. (n.d.). Classification: ROC curve and AUC. *Google for Developers*. <https://developers.google.com/machine-learning/crashcourse/classification/roc-and-auc>