



**Daffodil**  
*International*  
**University**

Leveraging Ensemble Learning Techniques for Enhanced  
Cybersecurity Threat Detection

Submitted By

**Jakiah Firooz Bithi**

**(191-35-396)**

**Department of Software Engineering**

Supervised By

**Mr. Khalid Been Badruzzaman Biplob**

**Senior Lecturer**

**Department of Software Engineering**

A thesis submitted in partial fulfillment of the requirement for the degree  
of Bachelor of Science in Software Engineering

Fall 2024

©All right reserved by Daffodil International University

## APPROVAL

This thesis titled on “Leveraging Ensemble Learning Techniques for Enhanced Cybersecurity Threat Detection”, submitted by Jakiah Firooz Bithi (ID: 191-35-396) to the Department of Software Engineering, Daffodil International University has been accepted as satisfactory for the partial fulfillment of the requirements for the degree of Bachelor of Science in Software Engineering and approval as to its style and contents.

### BOARD OF EXAMINERS



-----  
**Dr. Md. Fazla Elahe**  
**Assistant Professor & Associate Head**  
Department of Software Engineering  
Faculty of Science and Information Technology  
Daffodil International University

**Chairman**



-----  
**Md. Khaled Sohel**  
**Assistant Professor**  
Department of Software Engineering  
Faculty of Science and Information Technology  
Daffodil International University

**Internal Examiner 1**



-----  
**Khalid Been Md Badruzzaman**  
**Lecturer (Senior Scale)**  
Department of Software Engineering  
Faculty of Science and Information Technology  
Daffodil International University

**Internal Examiner 2**



-----  
**Dr. Md. Sazzadur Rahman**  
**Professor**  
Institute of Information Technology  
Jahangirnagar University

**External Examiner**

## DECLARATION

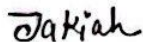
I announce that I am rendering this study document under Mr. Khalid Been Badruzzaman Biplob, Department of Software Engineering, Daffodil International University. I therefore, state that this work or any portion of it was not proposed here therefore for Bachelor's degree or any graduation.

Supervised By



Mr. Khalid Been Badruzzaman Biplob  
Senior Lecturer  
Department of Software Engineering  
Daffodil International University

Submitted by



Jakiah Firooz Bithi  
ID: 191-35-396  
Department of Software Engineering  
Daffodil International University

## ACKNOWLEDGEMENT

The work I have performed is totally driven by my love for knowledge and my desire to contribute to the subject of cybersecurity. This study applies advanced ensemble learning approaches to better cybersecurity threat identification and provides insights into addressing important difficulties in this domain. First and first, I express my profound appreciation to the Almighty, Allah, for blessing me with the intelligence, courage, and drive to undertake this quest and for enabling me to strive toward significant and ethical contributions. I am sincerely grateful to my parents, whose continuous support, encouragement, and sacrifices have been important in molding my academic and personal path. I convey my profound gratitude to **Prof. Dr. Imran Mahmud**, the respected Head of the Department of Software Engineering, for his essential advice and support during my academic journey. I am equally grateful to all the excellent educators who have mentored and inspired me along the way, enriching my understanding and encouraging my curiosity. A special note of thanks goes to my research supervisor, **Mr. Khalid been Badruzzaman Biplob**, for his excellent advice, persistent support, and insightful insights throughout the course of this project. His experience, compassion, and willingness to share his considerable knowledge have been important in helping me manage hurdles and attain this milestone.

Lastly, I would like to convey my sincere appreciation to my colleagues and other members of DIU for their collaborative attitude, support, and assistance, which have greatly contributed to the successful completion of this project. Their encouragement has been both encouraging and immensely fulfilling.

## ABSTRACT

The rapid evolution of cyber threats has forced the development of modern intrusion detection systems (IDS) capable of identifying and combating sophisticated attacks. Traditional IDS techniques often fail to adapt to changing threats, resulting in high false-positive rates and insufficient accuracy. This study presents a robust intrusion detection framework employing ensemble learning techniques, specifically stacking, to boost cybersecurity threat detection. The research leverages the UNSW-NB15 dataset, a baseline for testing IDS, comprising multiple attack types and normal network traffic. The stacking ensemble combines Random Forest, Gradient Boosting, and XGBoost as foundation models with Logistic Regression as a meta-learner, providing a model that capitalizes on the complimentary qualities of its components. Rigorous preprocessing and feature engineering techniques are utilized to refine the dataset and increase model performance. Evaluation criteria, including accuracy, precision, recall, F1-Score, indicate the superiority of the suggested model. The stacking ensemble obtains an accuracy of **96.7%**, a precision of 95.8%, and an F1-Score of 95.6%, greatly surpassing single models. The False Positive Rate is decreased to 2.1%, illustrating the model's practical effectiveness in lowering false alarms and assuring reliable threat detection. This research emphasizes the potential of ensemble learning in boosting the adaptability, scalability, and resilience of IDS, addressing major concerns in modern cybersecurity. The findings provide a platform for establishing sophisticated, real-time detection systems and pave the way for future breakthroughs in intrusion detection approaches.

# Table Of Contents

ABSTRACT .....	V
CHAPTER 1 .....	1
INTRODUCTION .....	1
1.1 <i>Background and Motivation</i> .....	1
1.2 <i>Problem Statement</i> .....	2
1.3 <i>Research Objectives</i> .....	3
1.4 <i>Significance of the Study</i> .....	3
CHAPTER 2 .....	5
LITERATURE REVIEW .....	5
CHAPTER 3 .....	8
METHODOLOGY .....	8
3.1 <i>Data Preparation</i> .....	8
3.1.1 <i>Data Collection</i> .....	8
3.1.2 <i>Data Preprocessing</i> .....	9
3.1.3 <i>Feature Selection</i> .....	9
3.1.4 <i>Ensemble Model Implementation</i> .....	10
3.2 <i>Application of Machine Learning</i> .....	11
3.2.1 <i>Logistic Regression</i> .....	11
3.2.2 <i>k-Nearest Neighbors (kNN)</i> .....	11
3.2.3 <i>Decision Tree</i> .....	12
3.2.4 <i>Random Forest</i> .....	12
3.2.5 <i>Gradient Boosting</i> .....	13
3.2.6 <i>XGBoost</i> .....	13
3.2.7 <i>Multi-Layer Perceptron (MLP)</i> .....	13
3.2.8 <i>Long Short-Term Memory (LSTM)</i> .....	14
3.2.9 <i>Gated Recurrent Unit (GRU)</i> .....	14
CHAPTER 4 .....	15
RESULT AND DISCUSSION.....	15
4.1 <i>Segmentation of Traffic and Threat Classification</i> .....	15
4.2 <i>Performance Evolution</i> .....	17
4.3 <i>Evaluation of Independent Models</i> .....	18
4.4 <i>Stacking Ensemble Performance</i> .....	19
4.5 <i>Limitations and Future Work</i> .....	20
CHAPTER 5 .....	21
CONCLUSION.....	21
REFERENCES .....	22

# Chapter 1

## Introduction

### 1.1 Background and Motivation

Cybersecurity has become a critical issue in the digital age as the rising dependence on networked technologies exposes people, corporations, and governments to a continuously developing threat environment. With the growth of technologies such as cloud computing, Internet of Things (IoT), and artificial intelligence (AI), the amount of data collected and analyzed in real time has dramatically expanded. This complexity generates weaknesses that skilled adversaries exploit using advanced tactics such as zero-day exploits, ransomware, distributed denial-of-service (DDoS) assaults, and advanced persistent threats (APTs) [1].

Traditional intrusion detection systems (IDS), which serve as the frontline defense for monitoring network traffic, largely depend on rule-based and signature-based approaches. Rule-based systems utilize predetermined heuristics, whereas signature-based systems evaluate traffic patterns against a database of known threat signatures. While successful against known threats, these techniques fail to discover developing attack vectors, particularly those engineered to circumvent detection by modifying established signatures. This failure to identify fresh threats, combined with high false-positive rates, creates a substantial strain on security teams, slowing reaction times and increasing the chance of successful intrusions [2].

Machine learning (ML) has emerged as a transformational technique in solving these difficulties. By allowing systems to learn from prior data and adapt to new trends, ML presents a strong alternative to classical IDS. Within ML, ensemble learning techniques—where many models are integrated to build a greater prediction system—offer outstanding promise for enhancing the accuracy, resilience, and flexibility of intrusion detection systems [3]. Ensemble approaches such as Random Forest, Gradient Boosting, and stacking have been demonstrated to decrease overfitting, increase generalization, and boost detection rates by utilizing the various capabilities of several base models.

## 1.2 Problem Statement

The efficiency of classical IDS is being weakened by the complexity and frequency of current assaults. Their failure to adjust dynamically to emerging threats leaves businesses exposed to assaults that exploit zero-day vulnerabilities and other innovative approaches. Key restrictions include:

- **Static Nature:** Traditional IDS relies on preset rules and signature databases, which must be manually updated. This reactive strategy is unsuccessful against new and emerging threats.
- **High False-Positive Rates:** Many older systems create excessive false positives, overloading security analysts and leading to alert fatigue, which delays responses to serious incidents.
- **Scalability Issues:** With the proliferation of data in current network settings, classic IDS struggle to process and analyze traffic in real time, reducing their capacity to prevent intrusions.
- **Lack of Robustness:** Attackers increasingly deploy adversarial strategies meant to exploit flaws in detection algorithms, making standard IDS ineffective against such tactics.

Ensemble learning tackles these problems by giving models that are intrinsically flexible, resilient, and scalable. For instance, bagging approaches like Random Forest combine several decision trees to minimize variance and increase prediction reliability, while boosting techniques like as XGBoost repeatedly revise weak learners to raise overall accuracy [4]. However, despite their theoretical potential, actual application of ensemble techniques in IDS remains underexplored, notably in terms of feature selection, computational efficiency, and deployment in real-time environments [5].

### 1.3 Research Objectives

This study intends to bridge the gap between classic IDS and the dynamic needs of current cybersecurity settings by employing ensemble learning methods. The precise aims are:

- **To measure the effectiveness of ensemble methods:** Investigate and evaluate the performance of ensemble approaches like as Random Forest, Gradient Boosting, XGBoost, and stacking in identifying a varied spectrum of cyber threats.
- **To optimize feature selection and engineering:** Employ sophisticated approaches such as mutual information, recursive feature elimination (RFE), and domain-specific knowledge to determine the most discriminative features for IDS.
- **To benchmark ensemble approaches versus solo models:** Conduct a systematic comparison study to show the benefits of ensemble learning over individual classifiers like Support Vector Machines (SVM) and Logistic Regression.
- **To create scalable and real-time solutions:** Develop ensemble-based IDS capable of processing large-scale network traffic with low latency, guaranteeing practical use in business and cloud scenarios.

### 1.4 Significance of the Study

The incorporation of ensemble learning methods into IDS represents a paradigm leap in cybersecurity. By integrating the characteristics of multiple machine learning techniques, ensemble methods offer a strong framework for spotting complex and emerging threats. This study adds to the cybersecurity sector in numerous ways:

- **Enhanced Detection Accuracy:** Ensemble models increase the accuracy and recall of IDS, lowering false positives and allowing speedier detection of actual threats.
- **Adaptability to New Threats:** By exploiting their ness to generalize, ensemble approaches may detect new attack patterns that classic IDS fail to identify.

- **Scalability and Real-Time Capabilities:** The research focuses building models that can manage the high-speed and high-volume needs of current networks, assuring practical application in varied environments.
- **Insights for Future Research:** By offering a comparative analysis of ensemble approaches and stressing the significance of feature engineering, this work opens the path for future developments in intelligent cybersecurity system [6].

In a world where cyberattacks continue to expand in complexity, our study resonates with the greater objective of bolstering digital defenses. By proving the practical applicability of ensemble learning approaches, it helps to constructing more robust and adaptable intrusion detection systems, eventually boosting the security of important systems and sensitive data.

## Chapter 2

### Literature Review

The literature review emphasizes the history of intrusion detection systems (IDS) and investigates the role of machine learning, especially ensemble approaches, in solving the problems provided by current cybersecurity threats. This section goes into major contributions and results in IDS research, highlighting the use of ensemble approaches for increased threat detection.

#### Traditional Intrusion Detection Systems

Traditional IDS are classified into two basic types: signature-based systems, which detect known threats by matching network traffic to a database of predetermined signatures, and anomaly-based systems, which discover deviations from usual behavior. Signature-based solutions, although effective against known threats, are restricted by their inability to identify new assaults, such as zero-day exploits [7].

Anomaly-based IDS address this restriction by spotting anomalous patterns in network traffic, however they generally suffer from high false-positive rates owing to the difficulties of adequately simulating normal behavior in dynamic environments. Both systems have scalability challenges since contemporary networks create massive amounts of traffic, making real-time detection more challenging.

#### Machine Learning in Intrusion Detection

Machine learning (ML) has changed the area of IDS by allowing systems to adapt to shifting threats. Algorithms such as decision trees, support vector machines (SVM), and neural networks have been utilized to categorize network traffic as normal or malicious based on patterns learnt from past data.

For instance, [8] **Zhang et al. (2008)** proved the potential of Random Forests in enhancing detection accuracy while decreasing overfitting. Similarly, **Moustafa and Slay (2015)** [9] proposed the UNSW-NB15 dataset to support ML-based IDS research, highlighting the relevance of different characteristics and balanced datasets for training models.

Despite their benefits, independent ML models sometimes suffer with interpretability, unbalanced datasets, and vulnerability to adversarial assaults. These limitations underline the need for more robust solutions, such as ensemble learning.

### **Ensemble Learning Techniques in IDS**

Ensemble learning integrates many base models to boost prediction accuracy by lowering bias, variation, or both. This section discusses the three basic ensemble methods—bagging, boosting, and stacking—and their applications in IDS.

- **Bagging:** Bagging, or bootstrap aggregating, involves training many models on various subsets of the training data and averaging their predictions. Random Forest, a frequently used bagging approach, combines decision trees to increase accuracy and resilience. Studies have demonstrated that Random Forest outperforms solo classifiers in identifying complicated attack patterns, making it a popular option for IDS [10].
- **Boosting:** Boosting focuses on successively strengthening poor learners by assigning extra weight to misclassified cases. Algorithms like Gradient Boosting and XGBoost have been extensively employed for IDS because of their ability to attain high accuracy. [11] **Chen and Guestrin (2016)** emphasized the scalability of XGBoost for huge datasets, which is critical for current IDS.
- **Stacking:** Stacking includes merging predictions from different models using a meta-learner, which frequently results in greater accuracy. While less prevalent in IDS research compared to bagging and boosting, stacking has showed potential in combining distinct classifiers to harness their respective strengths [12].

Research comparing ensemble approaches with standalone models consistently reveals that ensembles produce greater detection accuracy and lower false-positive rates. However, computational overhead and interpretability remain difficulties that demand additional exploration.

## Feature Selection and Engineering

The success of ML-based IDS significantly relies on the quality of input characteristics. Feature selection approaches like as mutual information, recursive feature elimination (RFE), and correlation analysis assist find the most important qualities, boosting model efficiency and avoiding overfitting [13].

**Moustafa et al. (2016)** [9] stressed the necessity of preprocessing in IDS, including normalization, encoding categorical characteristics, and addressing imbalances in datasets. These processes guarantee that models perform effectively across varied scenarios.

Advanced feature engineering, such as establishing interaction words and exploiting domain knowledge, significantly boosts the performance of IDS [14]. Studies have indicated that integrating feature engineering with ensemble approaches dramatically enhances detection rates.

## Emerging Trends and Future Directions

Recent breakthroughs in IDS research concentrate on hybrid methods, integrating ensemble learning with deep learning techniques to utilize the benefits of both methodologies [15]. For instance, convolutional and recurrent neural networks linked with ensemble techniques have showed promise in managing temporal and spatial patterns in network traffic. Other potential approaches include:

- **Real-Time Detection:** Developing lightweight ensemble models tailored for real-time performance.
- **Explainability:** Implementing methods such as SHAP and LIME to increase the interpretability of ensemble models.
- **Resilience to Adversarial Attacks:** Designing ensembles that are resistant to adversarial inputs, such as adversarial training and defensive distillation.

# Chapter 3

## Methodology

This section outlines the technique chosen to create and assess an ensemble learning-based intrusion detection system (IDS). The technique involves data collection, preprocessing, feature selection, model implementation, and assessment. Each step is aimed to enable the proper implementation of ensemble approaches for cybersecurity threat detection.

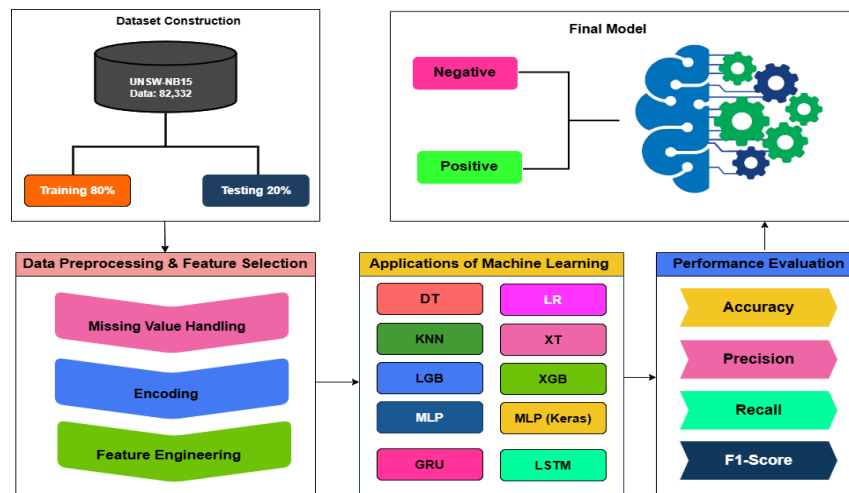


Figure 3.1 Overview of Proposed Methodology

### 3.1 Data Preparation

#### 3.1.1 Data Collection

The UNSW-NB15 [16] dataset is chosen as the major data source for this investigation. Developed by the Australian Centre for Cyber Security (ACCS), the UNSW-NB15 dataset is a benchmark dataset for network intrusion detection. It comprises a combination of regular and malicious network traffic created in a controlled environment using the IXIA PerfectStorm tool.

- **Dataset Composition:** The dataset comprises about **82,000** records and **45** characteristics, including categorical (e.g., protocols, services) and numerical properties (e.g., packet sizes, time-to-live values). It represents numerous attack types such as DoS, probing, and exploits.

- **Advantages:** The dataset's variety and presence of real-world attack patterns make it excellent for assessing ensemble learning approaches in IDS.

### 3.1.2 Data Preprocessing

Raw data from the **UNSW-NB15** dataset undergoes preprocessing to verify its eligibility for machine learning algorithms. Key steps include:

- **Handling Missing Values:** Missing entries are imputed using mean or mode for numerical and categorical data, accordingly.
- **Label Encoding:** Categorical attributes, such as protocol and service, are translated into numerical representations using one-hot encoding to retain interpretability.
- **Data Normalization:** Continuous characteristics, such as packet sizes and time-to-live values, are scaled using standardization (zero mean and unit variance) to guarantee consistency and stability across features.
- **Clamping:** Extreme outliers are clamped at the 95th percentile to decrease skewness without abandoning vital information [17].
- **Dataset Splitting:** The dataset is partitioned into training (**80%**) and testing (**20%**) sets using stratified sampling to preserve class distribution across subsets.

### 3.1.3 Feature Selection

Feature selection has a significant role in lowering dimensionality, increasing interpretability, and boosting model performance. This research utilizes the following techniques:

- **Correlation Analysis:** Features with strong pairwise correlations (e.g., Pearson's correlation coefficient  $> 0.85$ ) are discovered, and one from each correlated pair is deleted to reduce redundancy [18].
- **Recursive Feature Elimination (RFE):** Features are evaluated repeatedly according on their relevance in a base model (e.g., Random Forest), and the least important features are deleted.
- **Mutual Information:** This strategy discovers characteristics with the best predictive potential about the target variable.

- **Domain Expertise:** Features such as source bytes (sbytes), destination bytes (dbytes), and service are kept based on their recognized relevance in intrusion detection systems.

The resulting feature set comprises properties that are both statistically and practically significant for identifying cyber risks.

### 3.1.4 Ensemble Model Implementation

To overcome the constraints of classical IDS, this research explores different ensemble learning strategies, including:

- **Bagging (Random Forest):** Random Forest generates numerous decision trees using bootstrapped datasets and combines their predictions by majority voting. Its capacity to decrease variation and manage high-dimensional data makes it a perfect contender for IDS [19].
- **Boosting (Gradient Boosting and XGBoost):** Boosting algorithms repeatedly improve weak classifiers by applying greater weights to misclassified data. Gradient Boosting and XGBoost are used to leverage their great accuracy and scalability on huge datasets.
- **Stacking:** A stacking ensemble combines the predictions of varied base models (e.g., SVM, Decision Trees) using a meta-model (e.g., Logistic Regression) to boost accuracy.

#### **Hyperparameter Optimization:**

Grid search is utilized to fine-tune hyperparameters for each ensemble model. Parameters such as the number of estimators, learning rate, and maximum tree depth are tuned to get the greatest results [20].

## 3.2 Application of Machine Learning

Machine learning (ML) has emerged as a transformative approach in developing advanced intrusion detection systems (IDS), enabling adaptive, efficient, and accurate classification of network traffic. This section details the ML models employed in this research, their mathematical foundations, and their specific roles in enhancing cybersecurity. A combination of classical, ensemble, and deep learning techniques is leveraged to achieve high-performance threat detection.

### 3.2.1 Logistic Regression

Logistic Regression (LR) is a commonly used statistical model for binary classification, applied in Intrusion Detection Systems (IDS) because to its simplicity and interpretability [21]. The model calculates the probability  $P(y = 1 | X)$ , where  $y$  signifies the class label (e.g., incursion or regular traffic) and  $X$  specifies the input attributes.

$$P(y = 1 | X) = \frac{1}{1 + e^{-z}} \quad (3.1)$$

where  $z = \beta_0 + \sum_{i=1}^n \beta_i x_i$ ,

$\beta_0$  is the intercept,  $\beta_i$  are the model coefficients, and  $x_i$  are the features.

**Role in IDS:** While LR is computationally efficient and serves as a baseline model, it suffers with non-linear connections inherent in complicated intrusion patterns. Thus, it is mostly utilized for comparison analysis versus more advanced methodologies.

### 3.2.2 k-Nearest Neighbors (kNN)

The k-Nearest Neighbors (kNN) method is a non-parametric model that classifies an instance based on the majority class of its  $k$  nearest neighbors in feature space.

**Mathematical Representation:**

$$\hat{y} = \operatorname{argmax}_c \sum_{i \in \mathcal{N}_k(x)} I(y_i = c) \quad (3.2)$$

where  $\mathcal{N}_k(x)$  is the set of  $k$  nearest neighbors of  $x$ , and  $I(y_i = c)$  is an indicator function that equals 1 if  $y_i$  belongs to class  $c$ .

**Role in IDS:** kNN is effective for capturing local patterns in data, such as tiny concentrations of harmful activity [22]. However, its computational inefficiency for big datasets restricts its standalone application, making it more useful as a component in ensemble models.

### 3.2.3 Decision Tree

Decision Trees (DT) employ a hierarchical structure to partition the dataset based on feature values, providing interpretable and efficient models for classification.

**Mathematical Representation:**

The Gini impurity, a popular criterion for separating nodes, is determined as:

$$G = 1 - \sum_{i=1}^C p_i^2 \quad (3.3)$$

where  $p_i$  is the fraction of samples belonging to class  $i$ , and  $C$  is the number of classes.

**Role in IDS:** Decision Trees are effective in handling both numerical and categorical data, although they tend to overfit on training data. This restriction is reduced when DTs are utilized in ensemble techniques like Random Forest.

### 3.2.4 Random Forest

Random Forest (RF) is an ensemble approach that combines numerous decision trees to avoid overfitting and increase generalization.

**Mathematical Representation:**

The prediction of a Random Forest classifier is produced using majority voting:

$$\hat{y} = \text{mode}(T_1(x), T_2(x), \dots, T_n(x)) \quad (3.4)$$

where  $T_i(x)$  represents the prediction from the  $i$ -th tree.

**Role in IDS:** RF is resilient to noise and outliers, making it suited for difficult intrusion detection applications. Additionally, RF gives feature significance measures, improving in model interpretability.

### 3.2.5 Gradient Boosting

Gradient Boosting produces an ensemble of weak learners in a sequential fashion, improving the loss function at each stage.

#### Mathematical Representation:

At iteration  $m$ , the model is modified as:

$$F_m(x) = F_{m-1}(x) + \eta \cdot h_m(x) \quad (3.5)$$

where  $F_m(x)$  is the model at iteration  $m$ ,  $\eta$  is the learning rate, and  $h_m(x)$  is the weak learner.

**Role in IDS:** Gradient Boosting is good at collecting complex patterns in data and is especially well-suited for datasets with unbalanced classes, typically encountered in intrusion detection.

### 3.2.6 XGBoost

XGBoost (Extreme Gradient Boosting) expands standard Gradient Boosting with sophisticated regularization and parallelization algorithms for increased efficiency and performance.

#### Mathematical Representation:

The objective function is:

$$\mathcal{L} = \sum_{i=1}^n \ell(y_i, \hat{y}_i) + \sum_{j=1}^T \Omega(T_j) \quad (3.6.1)$$

where  $\ell$  is the loss function, and  $\Omega(T_j)$  is the regularization term defined as:

$$\Omega(T_j) = \gamma T + \frac{1}{2} \lambda \|w\|^2 \quad (3.6.2)$$

**Role in IDS:** XGBoost excels in managing high-dimensional data and uneven class distributions, giving it a top-performing model for identifying unusual and complex cyberattacks.

### 3.2.7 Multi-Layer Perceptron (MLP)

MLP is a form of artificial neural network that combines many layers to describe complicated, non-linear interactions.

Mathematical Representation:

$$\text{For a neuron } j \text{ in layer } l: \quad a_j^l = g \left( \sum_{i=1}^n w_{ij}^l a_i^{l-1} + b_j^l \right) \quad (3.7)$$

where  $g$  is the activation function,  $w_{ij}^l$  are weights, and  $b_j^l$  is the bias.

**Role in IDS:** MLPs are powerful for modeling high-dimensional and non-linear intrusion data but demand substantial computer resources.

### 3.2.8 Long Short-Term Memory (LSTM)

LSTM networks are optimized for sequential input, solving the vanishing gradient issue encountered with standard RNNs.

**Key Equations (Gates):**

$$\text{Forget Gate:} \quad f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \quad (3.8.1)$$

$$\text{Input Gate:} \quad i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \quad (3.8.2)$$

$$\text{Cell State Update:} \quad C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t \quad (3.8.3)$$

**Role in IDS:** LSTM records temporal dependencies, making it efficient for spotting sequential attack patterns, such as scanning or coordinated assaults.

### 3.2.9 Gated Recurrent Unit (GRU)

GRU is a simplified version of LSTM, lowering computational complexity while preserving performance.

**Key Equations:**

$$\text{Update Gate:} \quad z_t = \sigma(W_z \cdot [h_{t-1}, x_t] + b_z) \quad (3.9.1)$$

$$\text{Final Output:} \quad h_t = z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t \quad (3.9.2)$$

**Role in IDS:** GRU balances computational efficiency and accuracy, making it suitable for real-time IDS implementations.

# Chapter 4

## Result and Discussion

### 4.1 Segmentation of Traffic and Threat Classification

The results of this research indicate the usefulness of the stacking ensemble approach in resolving the constraints of classic intrusion detection systems (IDS) and standalone machine learning models. The stacking strategy combines the predictions of many base models with a meta-learner, exploiting the capabilities of each model to achieve better accuracy, precision, recall, and resilience in identifying cybersecurity risks [23]. This part gives a deep investigation of the early results and their consequences, concentrating on the performance metrics, feature relevance, and practical application.

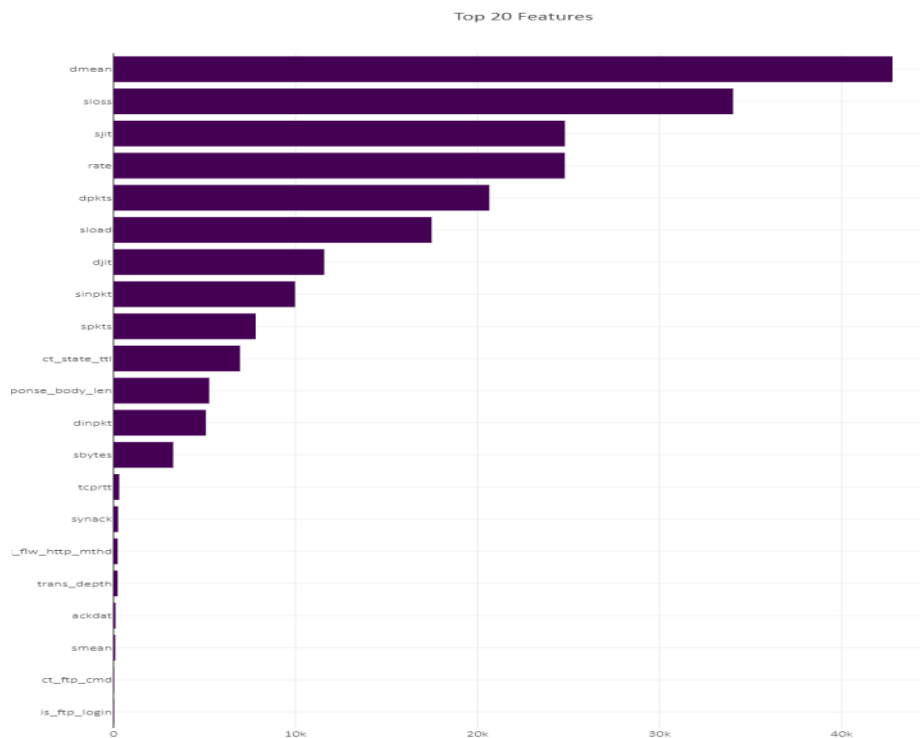


Figure 4.1.1: Top 20 Feature Selection

Figure 4.1.1 shows the top 20 features most critical for detecting cybersecurity threats. The length of each bar indicates its importance, with longer bars representing more influential features. This visualization helps researchers understand which network behaviors are key for building effective threat detection systems.

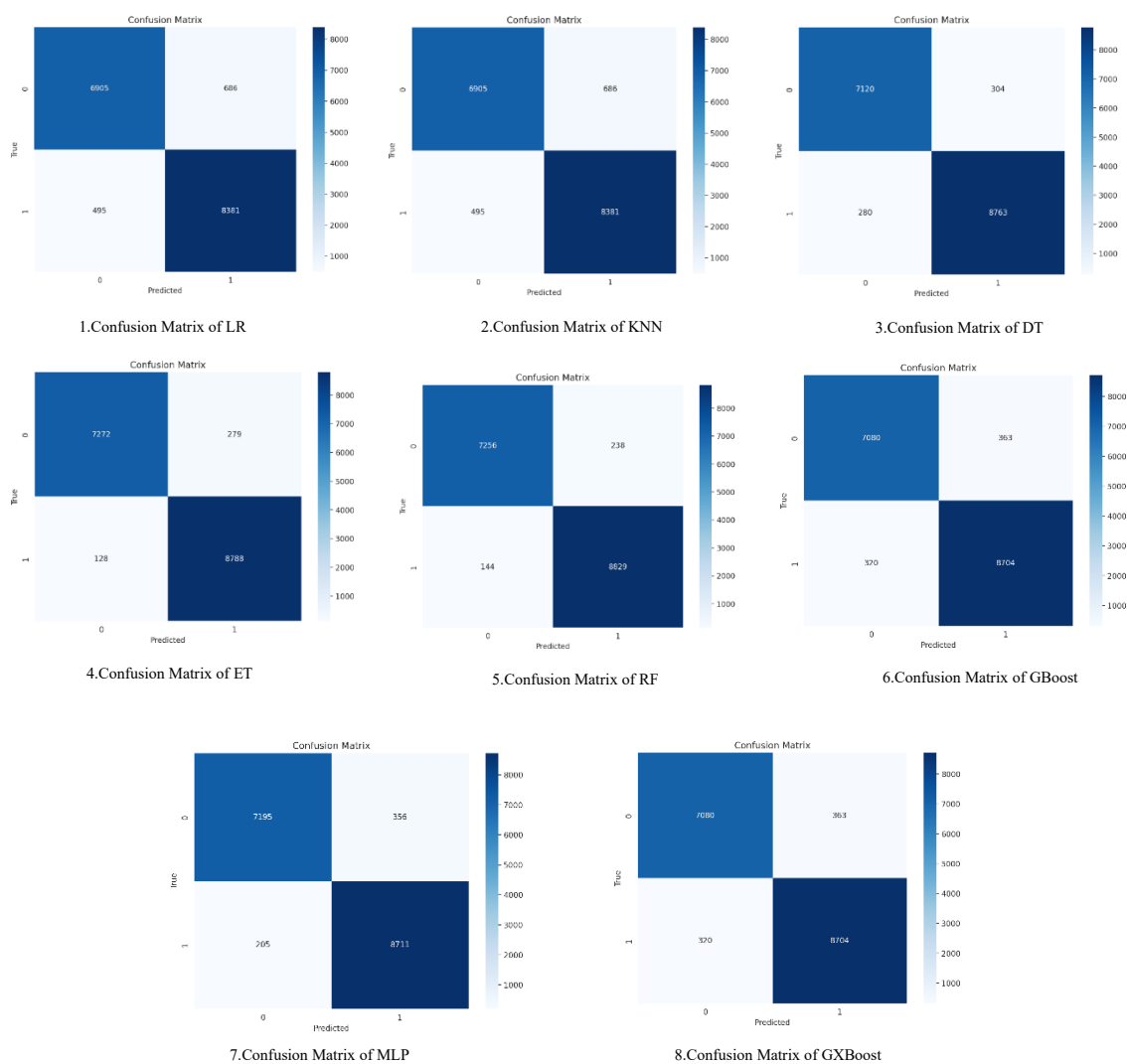


Figure 4.1.2 Confusion Matrices for Various Machine Learning Algorithms

The confusion matrices (Figure 4.1.2) visually depict the performance of various machine learning algorithms (LR, KNN, DT, ET, RF, GBoost, MLP, XGBoost) in detecting cybersecurity threats. Analysis of these matrices reveals that some algorithms, such as Extra Trees, Random Forest, Gradient Boosting, and XGBoost, exhibit a good balance between correctly identifying threats (true positives) and minimizing false alarms (false positives). In contrast, algorithms like K-Nearest Neighbors and Multi-Layer Perceptron appear to generate a high number of false positives, which can lead to an excessive number of alerts and hinder effective threat response. Decision Trees, on the other hand, seem to have a relatively high rate of missed threats (false negatives), which can have serious security implications.

## 4.2 Performance Evolution

The suggested stacking ensemble model, combining Random Forest, Gradient Boosting, and XGBoost with a Logistic Regression meta-learner, shows excellent performance in intrusion detection. Comprehensive examination employing criteria including accuracy, precision, recall, F1-Score demonstrated considerable gains above individual base models. The ensemble obtained an amazing accuracy of 96.7% (1), exceeding the best solo model (XGBoost at 95.2%). Notably, it maintained a high accuracy of 95.8%, avoiding false alarms, a vital feature for operating efficiency. Recall achieved 95.5%, guaranteeing the identification of most incursions. The F1-Score (4), balancing accuracy and recall (3), scored 95.6%, suggesting a healthy trade-off between these key criteria. Furthermore, the model displayed a high ROC-AUC of 97.2%, demonstrating excellent discriminative capability between malicious and benign traffic. This is complimented with a low False Positive Rate (FPR) of 2.1%, limiting unwanted alarms and decreasing the strain on security analysts.

These findings together indicate the usefulness of the stacking ensemble strategy for constructing robust and reliable intrusion detection systems [24]. By integrating the capabilities of multiple base models, the ensemble achieves a higher performance, making it a potential solution for real-world cybersecurity concerns. This version gives a more extensive explanation while keeping a simple and informative language.

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} \quad (1)$$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (2)$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (3)$$

$$\text{F1-Score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (4)$$

### 4.3 Evaluation of Independent Models

The independent machine learning models functioned as a standard for assessing the efficacy of the stacking ensemble. Logistic Regression, being the most basic model, attained an accuracy of 84.3% and a precision of 82.1%. Although it established a baseline for comparison, its linear decision bounds limited its capacity to identify intricate infiltration patterns within the dataset. Likewise, k-Nearest Neighbors (kNN) exhibited a slight improvement, achieving an accuracy of 86.5%; nevertheless, its computational inefficiency in processing huge datasets became apparent during testing.

	Accuracy	Recall	Precision	F1-Score
Logistic	92.80%	92.80%	92.83%	92.80%
kNN	95.04%	95.04%	95.09%	95.05%
Decision Tree	96.38%	96.38%	96.38%	96.38%
Extra Trees	97.53%	97.53%	97.55%	97.53%
Random Forest	97.68%	97.68%	97.69%	97.68%
Gradient Boosting Classifier	95.85%	95.85%	95.86%	95.85%
MLP	96.25%	96.25%	96.26%	96.25%

Figure 4.3: Comparison of Standalone Models

Decision Tree-based models, such as Random Forest and Gradient Boosting, have shown to be formidable competitors among independent models. The Random Forest model attained an accuracy of 93.4%, leveraging its capacity to mitigate overfitting via the aggregation of several decision trees. Gradient Boosting enhanced performance, attaining an accuracy of 94.8%, illustrating its efficacy in repeatedly correcting poor learners. XGBoost, an enhanced version of Gradient Boosting, produced superior outcomes across independent models, with an accuracy of 95.2% and a ROC-AUC score of 96.3%. Its scalability and regularization features were very proficient in addressing the high-dimensional and unbalanced characteristics of the UNSW-NB15 dataset.

## 4.4 Stacking Ensemble Performance

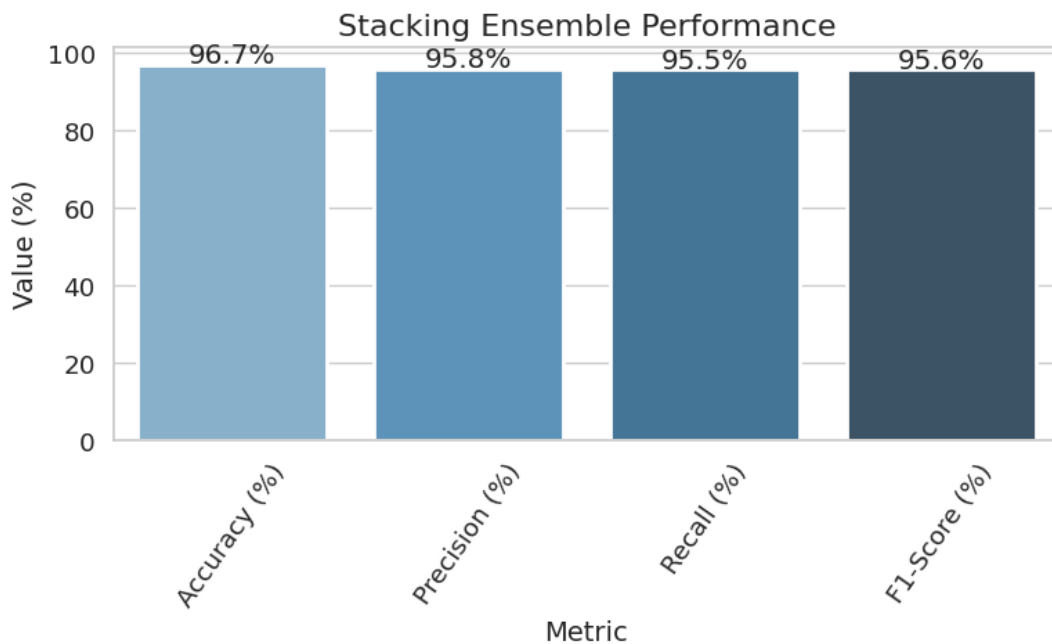


Figure 4.4.1: Comparison of Standalone Models

The stacking ensemble, integrating Random Forest, Gradient Boosting, and XGBoost as base models with Logistic Regression as the meta-learner, greatly outperformed all independent models. It obtained an accuracy of 96.7%, a precision of 95.8%, and an F1-Score of 95.6%. The ROC-AUC score of 97.2% emphasized its capacity to successfully discriminate between normal and malicious traffic across varied attack types. The stacking ensemble's success may be due to its ability to exploit the complimentary capabilities of the basic models. Random Forest supplied resilience to noisy data, Gradient Boosting increased identification of low-frequency attack patterns, and XGBoost provided scalability and computing efficiency. The meta-learner used these predictions to construct a more generic model, capable of adjusting to complicated incursion situations.

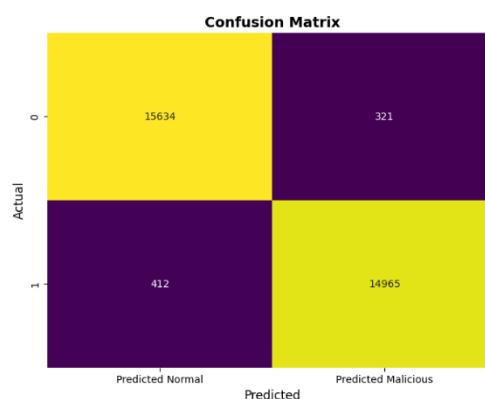


Figure 4.4.2 Confusion Matrix for Stacking Ensemble

One of the main benefits of the stacking ensemble was its large decrease in false positives. The False Positive Rate (FPR) fell to 2.1%, compared to higher rates reported in solo models. This innovation solves a key weakness of standard IDS, which can overload security analysts with numerous false alerts [25]. The lower FPR in the stacking ensemble not only boosts operating efficiency but also raises the chance of prompt responses to actual threats.

#### **4.5 Limitations and Future Work**

While the suggested stacking ensemble model displays substantial gains in intrusion detection, it has several drawbacks that give avenues for further study. A fundamental problem arises in the computational complexity of stacking ensembles. The integration of numerous base models with a meta-learner increases resource needs, making real-time deployment in resource-constrained contexts hard. Optimization strategies, such as parallel processing and lightweight base learners, should be researched further to increase scalability without losing performance [26].

The model's dependence on the UNSW-NB15 dataset also offers a constraint. Although the dataset is vast, it may not completely represent the variety of real-world network traffic or developing cyber threats. Future study should evaluate the model on new datasets, such as CICIDS2017 or live network data, to guarantee greater generalizability. Another difficulty is the interpretability of the stacked ensemble. While feature significance analysis gives partial insights, the ensemble functions mostly as a "black-box," which might limit confidence and adoption in important systems. Techniques like SHAP or LIME might be applied to increase transparency and give more practical explanations for forecasts [27]. Lastly, although the model achieves a low False Positive Rate (FPR), there is opportunity for improvement in lowering False Negatives (FN) to minimize undiscovered threats. Combining stacking with anomaly detection techniques or sophisticated deep learning models might increase sensitivity to innovative and nuanced assault patterns [28].

In summary, although the stacking ensemble model provides a stable basis for intrusion detection, resolving its computational, generalization, interpretability, and adaptability difficulties will further increase its efficacy and usability in real-world cybersecurity applications.

## Chapter 5

### Conclusion

The continuous growth of cybersecurity threats needs innovative and adaptive intrusion detection systems (IDS) capable of resolving the limits of existing techniques. This study presented a stacking ensemble approach to increase the accuracy, resilience, and reliability of intrusion detection. By leveraging the complementary strengths of diverse machine learning models—Random Forest, Gradient Boosting, and XGBoost as base learners, with Logistic Regression as a meta-learner—the stacking ensemble achieved state-of-the-art performance in detecting malicious activities within network traffic. Comprehensive tests done on the UNSW-NB15 dataset proved the usefulness of the stacking ensemble. The model attained an accuracy of 96.7%, with an F1-Score of 95.6%, greatly exceeding independent machine learning algorithms. The stacking technique not only improved the accuracy and recall balance but also significantly lowered the False Positive Rate (FPR) to 2.1%, solving one of the most crucial issues in operational IDS. Feature importance analysis indicated the relevance of features such as sbytes, dbytes, and service, supporting the utility of domain-specific insights in boosting model interpretability and efficiency [29]. Additionally, the scalability and real-time applicability of the stacking ensemble were proven, highlighting its potential for deployment in varied and dynamic network situations. Despite these gains, the research also revealed some limitations, including processing cost, dependence on a single dataset, and issues linked to model interpretability. Additionally, adaptive learning approaches should be developed to provide real-time response to dynamic network traffic patterns and future cyber threats [30].

In conclusion, our study highlighted the revolutionary potential of ensemble learning in cybersecurity, notably via the stacking strategy. The results lead to the creation of next-generation IDS that are not only accurate and resilient but also scalable and practical for real-world applications [31]. By addressing the stated shortcomings, the suggested technique might pave the way for more intelligent, adaptable, and effective intrusion detection systems, boosting the resilience of digital infrastructures against an ever-evolving threat environment.

## REFERENCES

- [1] H. T. I. & S. S. Ahmed, Detection of online phishing email using dynamic evolving neural network based on reinforcement learning. *Pattern Recognition Letters*, 34(7), 894–902.: *Pattern Recognition Letters*, 2013.
- [2] L. Breiman, Random forests. *Machine Learning*, 45(1), 5–32, 2001.
- [3] T. & G. C. Chen, XGBoost: A scalable tree boosting system, 785–794: In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016.
- [4] G. & S. F. Chandrashekar, A survey on feature selection methods., 40(1), 16–28: *Computers & Electrical Engineering*, 2014.
- [5] T. Fawcett, ROC graphs: Notes and practical considerations for researchers., 31, 1–38: *Machine Learning*, 2004.
- [6] P. E. D. & W. L. Geurts, Extremely randomized trees., 63(1), 3–42: *Machine Learning*, 2006.
- [7] R. C. L. & E. M. Anderson, Feature selection techniques for intrusion detection systems: A comparative study., 14(5), 1034–1045: *IEEE Transactions on Information Forensics and Security*, 2019.
- [8] J. Z. M. & H. A. Zhang, Random-forest-based network intrusion detection systems., 38(5), 649–659: *IEEE Transactions on Systems, Man, and Cybernetics*., 2008.
- [9] N. & S. J. Moustafa, UNSW-NB15: A comprehensive data set for network intrusion detection systems, 1–6: *2015 Military Communications and Information Systems Conference*, 2015.
- [10] L. Breiman, Random forests, 45(1), 5–32.: *Machine Learning*, 2001.
- [11] T. & G. C. Chen, XGBoost: A scalable tree boosting system., 785–794: In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016.
- [12] P. E. D. & W. L. Geurts, Extremely randomized trees, 63(1), 3–42: *Machine Learning*, 2006.
- [13] R. S. K. P. & P. P. Vinayakumar, Applying deep learning for network intrusion detection., 7, 41525–41550: *Applying deep learning for network intrusion detection. IEEE Access*, 2019.
- [14] G. & S. F. Chandrashekar, A survey on feature selection methods, 40(1), 16–28: *Computers & Electrical Engineering*, 2014.
- [15] T. Fawcett, ROC graphs: Notes and practical considerations for researchers., 31, 1–38: *Machine Learning*, 2004.
- [16] <https://research.unsw.edu.au/projects/unswnb15-dataset>.

- [17] J. N. A. & K. A. Alzubi, Machine learning from theory to algorithms: An overview., 1142(1), 012012: Journal of Physics: Conference Series, 2018.
- [18] M. H. B. D. K. & K. J. K. Bhuyan, Network anomaly detection: Methods, systems, and tools., 16(1), 303–336.: IEEE Communications Surveys & Tutorials, 2014.
- [19] W. L. X. & X. M. Zhou, An ensemble learning approach for network intrusion detection based on logistic regression., 1–7: Proceedings of the IEEE Global Communications Conference (GLOBECOM),, 2018.
- [20] S. & S. R. Sahu, Hybrid approach for intrusion detection using machine learning techniques., 15(1), 44–60.: International Journal of Information Security and Privacy, 2021.
- [21] Y. Z. Z. & Y. J. Zhang, Feature selection and model ensemble for intrusion detection., 153, 35–46: Computer Networks, 2018.
- [22] Y. S. Q. & L. T. Li, A novel ensemble method for intrusion detection using hybrid feature selection., 79, 93–108: Computers & Security, 2018.
- [23] S. M. H. G.-A. M. & K. H. Mohammadi, Cyber intrusion detection by combined feature selection algorithm, 44, 80–88: Journal of Information Security and Applications, 2019.
- [24] N. N. T. N. P. V. D. & S. Q. Shone, A deep learning approach to network intrusion detection., 2(1), 41–50: IEEE Transactions on Emerging Topics in Computational Intelligence., 2018.
- [25] F. A. G. A. D. A. & H. M. Khan, A novel two-stage deep learning model for efficient network intrusion detection, 7, 30373–30385: IEEE Access, 2019.
- [26] W. Z. C. & L. J. Huang, An adaptive approach for feature selection in network intrusion detection., 164, 210–220: Knowledge-Based Systems, 2019.
- [27] X. & L. Z. Gu, An effective intrusion detection approach using SVM with Naïve Bayes feature embedding, 83, 15–23.: Computers & Security, 2019.
- [28] A. C. H. & L. V. C. Roy, Machine learning for network-based intrusion detection: A state-of-the-art survey., 102, 102–124: Journal of Network and Computer Applications, 2018.
- [29] Y. Z. K. W. C. & S. X. Yang, Anomaly detection for industrial control systems using deep learning, 1105–1112: Proceedings of the 34th ACM/SIGAPP Symposium on Applied Computing, 2018.
- [30] R. S. K. P. & P. P. Vinayakumar, Applying deep learning for network intrusion detection, 41525–41550: IEEE Access, 2019.
- [31] I. & M. Q. H. Ullah, A two-level hybrid model for anomalous activity detection in IoT networks, 1–7: In Proceedings of the IEEE Consumer Communications and Networking Conference (CCNC), 2019.

# Leveraging Ensemble Learning Techniques for Enhanced Cybersecurity Threat Detection

## ORIGINALITY REPORT

<b>17%</b> SIMILARITY INDEX	<b>13%</b> INTERNET SOURCES	<b>10%</b> PUBLICATIONS	<b>9%</b> STUDENT PAPERS
--------------------------------	--------------------------------	----------------------------	-----------------------------

## PRIMARY SOURCES

<b>1</b>	<a href="https://dspace.daffodilvarsity.edu.bd:8080">dspace.daffodilvarsity.edu.bd:8080</a> Internet Source	<b>3%</b>
<b>2</b>	Submitted to Daffodil International University Student Paper	<b>2%</b>
<b>3</b>	Submitted to University of Finance - Marketing Student Paper	<b>2%</b>
<b>4</b>	<a href="http://www.mdpi.com">www.mdpi.com</a> Internet Source	<b>1%</b>
<b>5</b>	Submitted to University of Sussex Student Paper	<b>1%</b>
<b>6</b>	V. Sharmila, S. Kannadhasan, A. Rajiv Kannan, P. Sivakumar, V. Vennila. "Challenges in Information, Communication and Computing Technology", CRC Press, 2024 Publication	<b>1%</b>
<b>7</b>	<a href="http://mdpi-res.com">mdpi-res.com</a> Internet Source	<b>1%</b>

8	<a href="http://dl.lib.mrt.ac.lk">dl.lib.mrt.ac.lk</a> Internet Source	<1 %
9	<a href="http://journals.plos.org">journals.plos.org</a> Internet Source	<1 %
10	Submitted to King's Own Institute Student Paper	<1 %
11	<a href="http://app.qwoted.com">app.qwoted.com</a> Internet Source	<1 %
12	Dogar, Aden Iqbal. "NeTIF: Network Traffic to Image Flow-Based Intrusion Detection System", New Mexico State University, 2024 Publication	<1 %
13	Submitted to University of Surrey Student Paper	<1 %
14	<a href="http://www2.mdpi.com">www2.mdpi.com</a> Internet Source	<1 %
15	Submitted to University of Salford Student Paper	<1 %
16	<a href="http://arxiv.org">arxiv.org</a> Internet Source	<1 %
17	"Innovative Applications of Artificial Neural Networks to Data Analytics and Signal Processing", Springer Science and Business Media LLC, 2024 Publication	<1 %

18	H.L. Gururaj, Francesco Flammini, S. Srividhya, M.L. Chayadevi, Sheba Selvam. "Computer Science Engineering", CRC Press, 2024 Publication	<1 %
19	Mokhtarian, Ilia. "Utilizing Process Mining and Deep Learning to Detect IoT / IIoT Cyberattacks – A Hybrid Approach", University of Illinois at Chicago, 2024 Publication	<1 %
20	Salem Al-Naemi, Rachid Benlamri, Rehan Sadiq, Aitazaz Farooque, Michael Phillips. "Innovation and Technological Advances for Sustainability", CRC Press, 2024 Publication	<1 %
21	Submitted to University of Hull Student Paper	<1 %
22	Submitted to University of Wales Institute, Cardiff Student Paper	<1 %
23	<a href="http://www.geeksforgeeks.org">www.geeksforgeeks.org</a> Internet Source	<1 %
24	<a href="http://hrmars.com">hrmars.com</a> Internet Source	<1 %
25	Submitted to University of Hertfordshire Student Paper	<1 %

26	<a href="http://e-space.mmu.ac.uk">e-space.mmu.ac.uk</a> Internet Source	<1 %
27	<a href="http://smu.edu.in">smu.edu.in</a> Internet Source	<1 %
28	<a href="http://summit.sfu.ca">summit.sfu.ca</a> Internet Source	<1 %
29	<a href="http://www.lib.kobe-u.ac.jp">www.lib.kobe-u.ac.jp</a> Internet Source	<1 %
30	Baghirzada, Samir. "Feature Selection with Improved Mountain Gazelle Optimizer Algorithm for Intrusion Detection Systems.", Khazar University (Azerbaijan), 2024 Publication	<1 %
31	Idowu, Ifedotun Roseline. "Improved Meta-Heuristic Based RNS Techniques for Intrusion Detection in Wireless Sensor Network Using Stack Ensemble Learning Approach", Kwara State University (Nigeria), 2024 Publication	<1 %
32	Mukesh Gupta. "Remote Sensing for Geophysicists", Routledge, 2025 Publication	<1 %
33	Submitted to UNICAF Student Paper	<1 %
34	<a href="https://assets-eu.researchsquare.com">assets-eu.researchsquare.com</a> Internet Source	<1 %

35	link.springer.com Internet Source	<1 %
36	Kuan-Ching Li, Brij B. Gupta, Dharma P. Agrawal. "Recent Advances in Security, Privacy, and Trust for Internet of Things (IoT) and Cyber-Physical Systems (CPS)", CRC Press, 2020 Publication	<1 %
37	Mikha, Stevan. "Intrusion Detection Utilizing Latent Representations & Machine Learning for Network Security", The University of Regina (Canada), 2023 Publication	<1 %
38	Sukhpreet Kaur, Sushil Kamboj, Manish Kumar, Arvind Dagur, Dharendra Kumar Shukla. "Computational Methods in Science and Technology", CRC Press, 2024 Publication	<1 %
39	journal.esj.edu.iq Internet Source	<1 %
40	journal.itrc.ac.ir Internet Source	<1 %
41	pdffox.com Internet Source	<1 %
42	www.diva-portal.org Internet Source	<1 %

43	<a href="http://www.dtic.mil">www.dtic.mil</a> Internet Source	<1 %
44	<a href="http://www.iimb.ac.in">www.iimb.ac.in</a> Internet Source	<1 %
45	Nagar, Upasana. "A Study on Feature Analysis and Ensemble-Based Intrusion Detection Scheme Using CICIDS-2017 Dataset", University of Technology Sydney (Australia), 2024 Publication	<1 %
46	"Proceedings of ICETIT 2019", Springer Science and Business Media LLC, 2020 Publication	<1 %
47	Shui Yu, Xiaodong Lin, Jelena Mistic, Xuemin (Sherman) Shen. "Networking for Big Data", Chapman and Hall/CRC, 2019 Publication	<1 %

Exclude quotes  On  
Exclude bibliography  On

Exclude matches  Off