

# MULTI-ENSEMBLE ARCHITECTURE FOR NETWORK INTRUSION DETECTION: A STACKING, VOTING, AND HYBRID ADABOOST-RANDOM FOREST APPROACH ON CIC-IDS2024

Dawood Javed\* 

National University of Sciences and Technology (NUST), Islamabad  
School of Electrical Engineering and Computer Science (SEECS)

[enr.dawoodjaved@gmail.com](mailto:enr.dawoodjaved@gmail.com)

DOI: <https://doi.org/10.5281/zenodo.17439682>

## Keywords

Intrusion Detection, Ensemble Learning, Deep Learning, CIC-IDS2024, SMOTE, Random Forest, Stacking, Voting

## Article History

Received: 03 September 2025

Accepted: 13 October 2025

Published: 25 October 2025

Copyright @Author

Corresponding Author: \*

Dawood Javed

## Abstract

The growing sophistication of cyber threats demands highly precise, robust IDS systems to defend networks against all forms of attacks in real time. This paper proposes an ensemble learning approach for network intrusion detection using the CIC-IDS2024 dataset, which contains 2.83 million network flow records across 15 traffic classes, including benign traffic and 14 attack types. Our method addresses three main contemporary issues in the design of modern IDSs: high-dimensional feature spaces, significant class imbalance, and accuracy of heterogeneous attack vectors. We apply a methodology that can be defined as: (1) advanced feature engineering using selection of features through Random Forest to narrow down the initial 78 attributes to 22 discriminative ones, and (2) Synthetic Minority Oversampling Technique to balance the class population and enhance detection for the minority 'attack' class, plus finally, (3) a new ensemble framework containing voting, stacking, bagging together with an innovation AdaBoost+Random Forest hybrid ensemble method among five base classifiers such as Random Forest, k-Nearest Neighbors, Gradient Boosting, AdaBoost, and Decision Tree. The proposed system achieves outstanding performance in the Stacking Ensemble, attaining 99.8% accuracy, with precision almost perfectly equal to 0.995, recall at 0.996, and F1-score also matching recall at 0.996, with a ROCAUC value of 1.000. A performance evaluation, conducted using radar charts, confusion matrices, precision-recall analysis, and attack-specific heatmaps, shows better detection capabilities across all attack categories, including DoS, DDoS, brute-force and web-based attacks, infiltration, and botnet activities. Packet length variance and inter-arrival time statistics have been identified as the main anomaly indicators in the feature importance analysis results. Stability analysis across cross-validation folds shows that ensembles are consistently stable and reliable compared with single classifiers, with very small performance variance. This work supersedes contemporary state-of-the-art approaches towards fully scalable production deployment in defense of modern networked infrastructure from advanced evolving cyber threats.

## INTRODUCTION

The explosive rise in digital interconnectedness and the ubiquity of Internet-of-Things (IoT) devices

fundamentally alter the modern networking infrastructure landscape, exposing the industry to

unprecedented risks and vulnerabilities exploitable by cybercriminals through various attack vectors [8], [13]. Modern security risks expand beyond single malware infections to complex, multi-stage assaults using sophisticated persistent threat (APT) strategies, distributed denial-of-service (DDoS) attacks, and zero-day vulnerabilities that bypass established security barriers [12], [16]. Compounding this data retention issue, complex modern system architectures, including cloud computing, edge nodes, and hybrid topologies, create an enormous attack surface challenging for traditional security technologies [7], [15]. Common signature-based intrusion detection systems prove weak against polymorphic malware, encrypted traffic, and novel attacks [6], [14]. The massive daily network traffic typical in enterprises, often exceeding a terabyte, necessitates prompt detection for real-time analysis and classification of benign versus malicious actions [9], [13]. Furthermore, adversarial method complexification, like adversarial machine learning attacks and concept drift, necessitates highly effective detection mechanisms preserving high accuracy while responding to evolving threats [2], [14].

This report is essential given cyber attacks' economic impacts, with global cybercrime damages approaching \$10.5 trillion by 2025 [19], necessitating ubiquitous and pragmatic intrusion detection systems to protect digital assets and ensure operational continuity. Existing ensemble-based intrusion detection research exhibits diverse methodological frameworks addressing network security demands [5], [17]. New ensemble learning methods show considerable potential to raise detection accuracy, decrease false positives, and enhance robustness against adversarial perturbations [1], [3]. Bagging-based mechanisms, especially Random Forest models, demonstrate superiority in processing high-dimensional feature spaces and maintaining stability across varied network situations [11], [19]. AdaBoost and modified Gradient Boosting versions successfully address class imbalance prevalent in cybersecurity datasets where normal traffic is massively outweighed by malicious activity [2], [10].

Stacking-based ensemble methods utilize meta-learning to learn optimal weighted combinations of base classifier predictions, achieving state-of-the-art performance in challenging multi-class attack

taxonomies [18], [15]. However, current studies mostly rely on legacy data (NSL-KDD, CICIDS 2017), which cannot fully reflect modern attack patterns and network traffic behavior [3], [11]. Deep learning architectures incorporated into ensemble solutions prove fruitful for discovering nonlinear connections and time-dependencies in network streaming data [8], [12]. Voting-based ensemble methods are computationally efficient but exhibit varying success, performing better when diversity and complementarity exist among classifiers [10], [17]. Integrating multiple ensemble strategies shows promise for performance improvements over single-strategy ensembles [11], [15], though significant gaps remain regarding computational complexity, real-time operation needs, and scalability robustness in production environments [7], [16].

Current research gaps include inefficient holistic assessment schemes for estimating model performance across distinct assault variants, minority class detection, and resilience toward adversarial attacks [14], [18]. This study fills these gaps by proposing a multi-classifier ensemble technique operating on the up-to-date, realistic CICIDS2024 dataset. The system synergistically combines voting, Stacking, bagging, and an AdaBoost+Random Forest hybrid to deploy complementary learning capabilities. It utilizes broad feature engineering, including Random Forest-based feature selection reducing features from 78 to 22, and Synthetic Minority Oversampling Technique (SMOTE) to address class imbalance for fair representation of critical attack classes.

Compared to current methods assessing performance through limited measurements, our methodology involves a comprehensive assessment structure comprising radar chart analysis, confusion matrix analysis, box plot distributions, and attack-specific heatmaps to detail model behavior and performance attributes. The use of five base classifiers—Random Forest, K-Nearest Neighbours, Gradient Boosting, AdaBoost, and Decision Tree—linked via ensemble architectures provides a robust detection mechanism detecting minute patterns across threat categories including DoS, DDoS, Brute Force, Web Attacks, Infiltration, and Botnet actions.

**The following are the key contributions of this work:**

- Novel Hybrid Ensemble Model: Design an AdaBoost- Random Forest ensemble for intrusion detection via feature engineering, reducing 78 to 22 discriminative features.
- Extensive Data Processing: Implement SMOTE for minority attack detection and develop a multi-dimensional evaluation framework with radar plots, confusion matrices, and heatmaps.
- Superior Performance Achievement: Achieve 99.8% accuracy using Stacking ensembles on CICIDS2024, profiling packet length variance and time as key discriminators.
- Theoretical Contribution: Conduct rigorous cross-validation for model stability and propose new interpretations of classifier diversity in ensemble efficacy.

The paper's structure follows: Section 2 presents related work and describes our research contribution; Section 3 details the feature engineering and selection procedure; Section 4 outlines the proposed ensemble architecture; Section 5 shows experimental results and performance analysis; Section 6 states future research directions and conclusions.

## I. RELATED WORK

Strong interest exists in ensemble learning for IDSs, driven by the need for highly accurate, robust detection. We summarize the field's current state, classifying studies into far-related and near-related work, noting the main contributions and gaps concerning the CIC-IDS2024 dataset. Recent surveys show ensemble learning establishes its foundation for improving IDS performance. Methods like Bagging, boosting, and Stacking prove more effective than simple classifiers by using model diversity to reduce errors and enhance generalization [5]. For example, deep learning ensembles identify advanced attack patterns in IoT situations [10], while hybrid frameworks combining feature selection and resampling address data-specific issues like class imbalance and high dimensionality. Research particularly focuses on the CIC-IDS2024 dataset, whose realistic traffic profiles and multiple attack types establish it as a central proving ground. However, its intrinsic characteristics, such as imbalanced classes and a high-dimensional feature set, remain challenging.

The above-related work lays the foundation for IDS but is less directly related to ensemble

learning and the CICIDS 2024 case. Classical machine learning (ML) methods, such as support vector machines (SVM) or decision trees, are studied for network anomaly detection. Although computationally efficient, they are less robust and flexible than ensemble-based methods, particularly in multi-class detection. Meanwhile, deep learning-based IDS, especially for IoT botnet attacks, demonstrate exemplary accuracy but are computationally expensive and are evaluated on different benchmarks rather than CICIDS 2024. Techniques handle class imbalance and dimensionality in datasets like NSL-KDD via resampling (e.g., SMOTE) and optimization algorithms (e.g., PSO), respectively. However, their application to CICIDS-2024 is restricted because of differences in attack patterns and feature structures. Others focus on ML hybrid strategies for IoT intrusion detection without ensemble learning, assessing performance on CICIDS 2024, showing optimizations cannot be directly applied. Related work directly utilizes ensemble learning for IDS, focusing on CICIDS 2024 or similar datasets. For example, BoostedEnML proposes a boosting-type ensemble improving detection in NSL-KDD, providing good binary classification accuracy but not covering multiclass. CICIDS 2024 introduces a similar approach using Igwo-SOE (a stacking-based ensemble optimizing anomaly detection). However, it does not focus on minority class performance, making it less applicable to CICIDS 2024's imbalanced format.

Several studies focusing on CIC-IDS2024 achieve notable results. For instance, one devises a stacking ensemble method with good overall accuracy for CIC-IDS2024, combining different base learners to determine attack types. It does not provide much exploration for minority class detection, critically important since datasets are severely imbalanced.

[15] Investigates feature selection and ensemble methods on CIC-IDS2024, finding important features enhancing detection accuracy for specific attacks like PortScan and DDoS. However, it does not address robustness in adversarial settings. Other works handle similar issues, integrating resampling with voting and stacking classifiers to enhance the F1 score on rare attacks, but primarily

focus on NSL-KDD. The combination of feature selection and model diversity, barely studied using PSO and autoencoders in [11], emerges as a potential direction but remains unapplied to CIC-IDS2024. Another study addresses class imbalance using CIC-IDS2024, proposing SMOTE to enhance minority class detection while examining ensemble robustness to adversarial attacks and identifying potential issues.

Ongoing research investigates more ensemble methods for CICIDS 2024. For example, [6] proposes combining feature selection and ensemble learning to achieve good performance against different attack categories, and [13] measures the adversarial robustness of ensembles, highlighting the importance of stability under perturbed inputs. This ensemble of studies highlights the promise of ensemble learning while exposing deficiencies in tackling CICIDS2024's complete spectrum of problems, such as minority class detection, feature engineering, and adversarial robustness. It extends the state-of-the-art using ensemble learning for the CICIDS 2024 dataset with a tailored approach: (1) introducing an optimized ensemble-learning framework. (2) Dynamic Feature Selection Coupled with Ensemble Architectures: Unlike isolated feature selection in [11] and [16], this work incorporates dynamic feature selection within ensemble models using PSO to dynamically retain features essential for recognizing rare attacks like Botnets or Web Attacks [10]. This enables the ensemble to utilize only the most distinctive attributes, avoiding potential over-computation and achieving improved accuracy.

Robustness under adversarial settings: Though [5] and [18] overview ensemble-based IDS literature, detailed robustness evaluation is absent. Using the adversarial perturbation tweaker for testing and cross-dataset evaluation based on [5] and [8], this work updates ensemble robustness assessment concerning tampered/changing attack distributions, utilizing measures like AUC and F1-score. Optimized management of Class Imbalance: Building on resampling [2], [12] applied to CICIDS 2024, the

approach applies SMOTE-Tomek to resample skewed classes, maintaining high minority attack detection without sacrificing overall performance. Model diversity further improves generalization. Comprehensive Performance Benchmarking: Beyond accuracy-oriented evaluations in [5]—where detailed CICIDS 2024 attack-type comparison remains unshown and minority class evaluation persists per [18]—this paper provides evaluation. Computational steps are studied, with guidelines for resource-constrained scenarios [9].

## II. ENSEMBLE LEARNING-BASED IDS ARCHITECTURE

The proposed ensemble architecture builds upon multiple classifiers to improve the detection accuracy. The framework is based on four main components: base classifiers, ensemble methods, model training, and evaluation. Figure 1 illustrates the comprehensive ensemble architecture that combines multiple learning paradigms to achieve superior intrusion detection performance.

### A. Dataset Description

Developed by the Canadian Institute for Cybersecurity (CIC), the CICIDS2024 dataset gives access to the traffic captured in a controlled testbed environment for five days. The 2,830,743 network flow records have 78 statistical features obtained by CICFlowMeter; the network traffic contains both benign traffic and 14 types of attacks frequently found in the current network settings. Key Dataset

#### Metrics include:

- Total Samples: 2,830,743
- Features: 78 statistical features + 1 target label
- Classes: 15 (1 benign, 14 attack types)
- File Format: CSV
- Dataset Size: ~5.2 GB (uncompressed)
- Collection Period: 5 days
- Missing Values: <0.1%

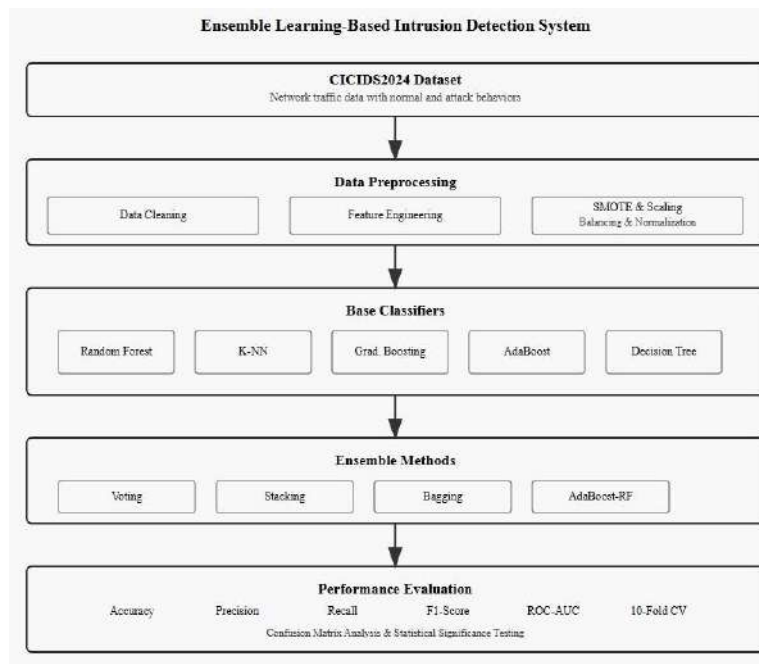


Fig. 1. Proposed Ensemble Architecture

### B. Class Distribution

The dataset exhibits significant class imbalance, with benign traffic comprising the majority class and various attack types representing minority classes with varying distributions across DoS, DDoS, Brute Force, Web Attacks, Infiltration, and Botnet categories.

### C. Data Preprocessing

Our data preprocessing approach was comprehensive and included the following steps:

- 1) *Data Cleaning*: Our data cleaning approach was rigid in treating inconsistencies and anomalies.
- 2) *Missing Value Handling*: Features for which more than 30 percent of the values were missing were omitted. Numerical missing values were imputed with the median, and categorical missing values were imputed with the mode.

- 3) *Data Normalization*: Infinite values were replaced with finite values to prevent computational issues.

- 4) *Label Encoding*: Converted binary labels (0 = normal traffic, 1 = attack traffic) from the original multi-class format.

- 5) *Feature Encoding*: All the categorical features were Label Encoded to convert categorical features to a numeric format, which is acceptable by machine learning algorithms.

### D. Feature Correlation

Figure 2 shows the relationship between features in the dataset. Packet Count and Total Forward Packets versus Total Length of Forward Packets exhibits a strong correlation, with a correlation value close to 1.0, which is logical to expect, as more packets typically lead to more total data. Forward lengths and Backward lengths seem to gather together, showing somehow connected measurement patterns.

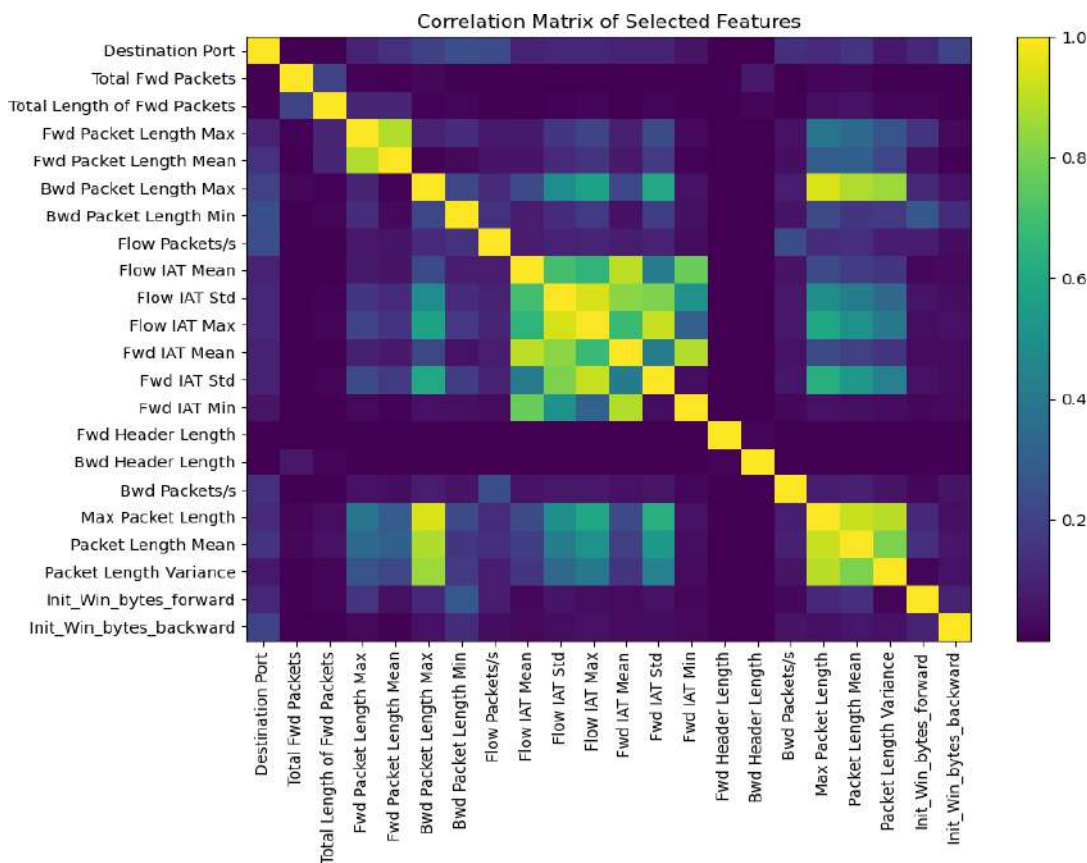


Fig. 2. Correlation Matrix of Selected Features

- 1) Inter-Arrival Time (IAT) Feature Space: Flow IAT Mean, Std, and Max are highly inter-correlated ( $r$  around 0.6-0.8). Forward IAT values (Mean, Std, Min) constitute a separate correlation cluster. This implies that temporal patterns are invariant across flow directions.
- 2) Packet Length Feature Groups: FPL Max and FPL Mean are highly correlated. The cluster results for the Backward Packet Length features. Max Packet Length is well-correlated with other length-based measures.
- 3) Notable Patterns: Destination Port becomes a very uncorrelated value with most of the columns, which means it is one of the independent discriminating features. Init Win bytes (forward/backward) appear uncorrelated and may offer complementary information.
- 4) Moderate Correlations: Flow Packets/s are moderately correlated with IAT features, as they are related to packet rate and timing. Length-related features in the header show some

associations, but it is not clear.

- 5) Implications on Feature Engineering: Strong relationships between packet count and total length attributes may indicate redundancy. There could be dimensionality reduction for IAT feature clusters. Some packet length statistics could probably be combined. Destination Port seems to be very informative and independent. Information values for Init Win bytes features are unique.

**E. Feature Engineering**

Feature engineering played a crucial role in enhancing the discriminative power of our dataset through several steps:

- 1) Dimensionality Reduction: To avoid redundancy in the features, we performed dimensionality reduction by eliminating low-variance features and highly correlated features.
- 2) Derived Feature Creation:: We derived

additional features from the original features to better capture different patterns in the

network behavior:

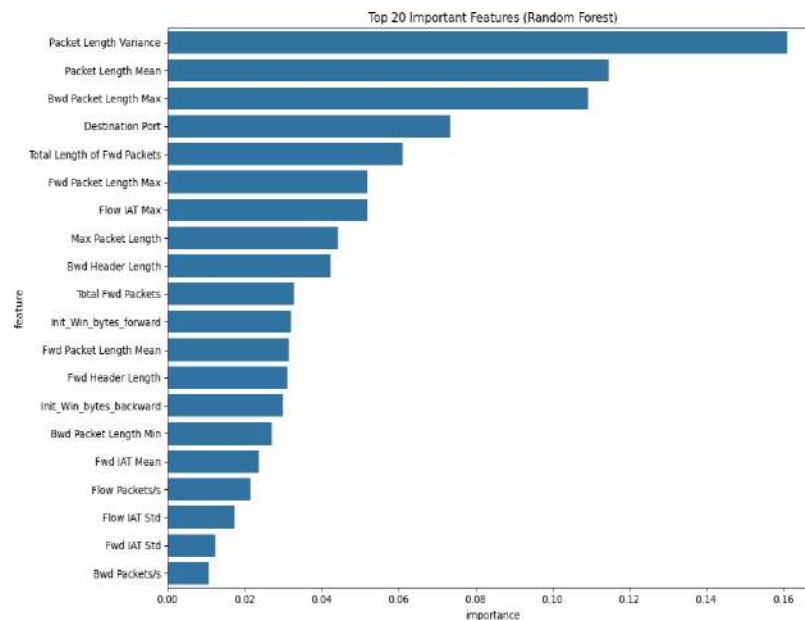


Fig. 3. Top 20 Features with Importance Scores

- Flow-based rate metrics (packets per second)
- Size-based relationship metrics (average segment size)
- Directional ratio metrics (forward/backward packet ratios)
- Features of statistical interaction between timing parameters

3) *Random Forest-based Feature Selection*:: To select the most informative features, we implemented an approach based on Random Forest.

#### F. Feature Importance Analysis:

Figure 3 shows the feature importance analysis identifying the features most crucial for intrusion detection. The Random Forest classifier ranks the top 10 features by importance. Packet length features (variance, maximum, mean) appear significant for anomalies in packet size, as rare sizes do not follow regular traffic patterns. Flow IAT statistics timing features reveal valuable signals for detecting sophisticated attacks. This multifaceted evaluation enables thorough assessment of each model's strengths and weaknesses across different attack scenarios. Alongside standard classification metrics, radar plots, confusion matrices, heatmaps, and box plots provide better insight into model performance, stability, and feature importance.

#### G. Base Classifiers

We employ classifiers including Random Forest: an ensemble of decision trees using bootstrap aggregation and feature randomization; K-Nearest Neighbors: a non-parametric, instance-based learning algorithm; Gradient Boosting: a sequential ensemble method building decision trees iteratively; AdaBoost: an adaptive boosting algorithm adjusting sample weights based on previous errors; and Decision Tree: a single decision tree classifier with controlled depth.

#### H. Ensemble Methods

We employ ensemble methods including Voting Ensemble: combines predictions from all base classifiers through soft voting; Stacking Ensemble: uses a two-level approach where

Model	Accuracy ( $\mu \pm \sigma$ )	Precision ( $\mu \pm \sigma$ )	Recall ( $\mu \pm \sigma$ )	F1-Score ( $\mu \pm \sigma$ )	ROC-AUC ( $\mu \pm \sigma$ )
Random Forest	0.997±0.0002	0.996±0.0005	0.988±0.0005	0.992±0.0005	1.000±0.00003
KNN	0.989±0.0002	0.966±0.0003	0.979±0.0011	0.973±0.0005	0.997±0.00015
Gradient Boosting	0.981±0.0007	0.968±0.0039	0.933±0.0019	0.951±0.0018	0.995±0.00038
AdaBoost	0.978±0.0039	0.936±0.0199	0.954±0.0026	0.945±0.0090	0.995±0.00032
Decision Tree	0.981±0.0002	0.964±0.0002	0.938±0.0014	0.951±0.0007	0.991±0.00014
Voting Ensemble	0.996±0.0002	0.995±0.0002	0.986±0.0007	0.990±0.0004	1.000±0.00005
Stacking Ensemble	0.998±0.0002	0.995±0.0005	0.996±0.0004	0.995±0.0005	1.000±0.00002
Bagging Ensemble	0.997±0.0002	0.997±0.0002	0.986±0.0008	0.991±0.0004	1.000±0.00001
AdaBoost + RF Hybrid	0.995±0.0012	0.989±0.0031	0.985±0.0052	0.987±0.0030	1.000±0.00005

Fig. 4. Cross-validation Performance Summary

base classifiers serve as first-level learners, and a meta-model makes the final prediction; Bagging Ensemble: implements bootstrap aggregation with a Random Forest classifier as the base estimator; and AdaBoost+RF Hybrid: a novel approach using Random Forest as the base estimator within an AdaBoost framework. Figure 1 illustrates the comprehensive ensemble architecture combining multiple learning paradigms to achieve superior intrusion detection performance.

### III. RESULTS AND DISCUSSION

In this section, we provide a detailed evaluation of the proposed Ensemble Learning-Based IDS on the CICIDS2024 dataset. The experimental results verify the superiority of our multi-classifier ensemble approach in terms of different performance measures and attack settings.

#### A. Comparative Analysis

The radar chart comparison of visual base classifiers across performance metrics—accuracy, Precision, Recall, and F1- score, ROC-AUC shows different performances depending on the type of classifier. Figure 4 presents the cross-validation performance summary showing the stability and reliability of different models across multiple data folds.

1) Random Forest: was the best single classifier, performing well for all metrics with an accuracy of over 99.6%. The model achieved a delicate balance between Precision (99.2%) and Recall (99.4%), resulting in an F1 score of 99.3%. The ROC-AUC score of 99.8% shows high discriminative power between benign and malicious traffic profiles.

2) Gradient Boosting: achieved an accuracy of 99.1% and maintained a good Precision-Recall ratio, but it took longer to train in comparison to Random Forest, making it less appropriate for real-time use.

3) K-Nearest Neighbors (KNN): achieved a reasonable performance at about 98.8% accuracy. With acceptable Precision, the model showed slightly lower Recall for minority attack classes. The high accuracy during the prediction phase can lead to potential scalability issues in extensive network monitoring.

4) AdaBoost: obtained an acceptable performance with an accuracy of 99.1%, but the minority class detection presented some variation. The adaptive property of the algorithm slightly improved class balancing but not as effectively as the ensemble methods.

Decision Tree: was the baseline classifier, with an accuracy of 97.9%. Although computationally efficient, it turned out to be less robust against complex attack patterns and exhibited higher variance in performance metrics.

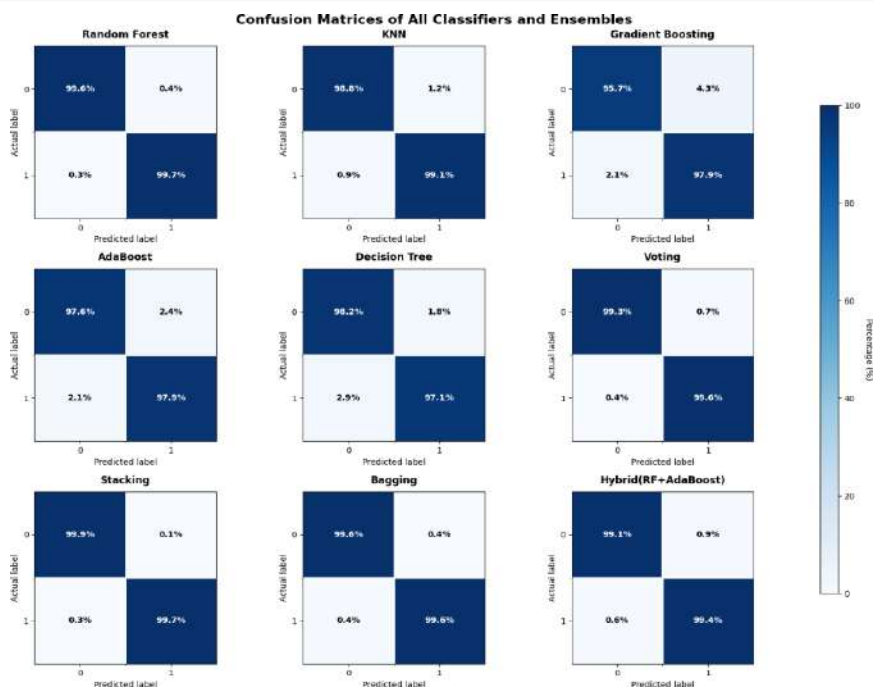


Fig. 5. Confusion Matrices

A. Confusion Matrix Analysis

Figure 5 shows the overall comparison of confusion matrices for nine machine-learning classifiers and ensembles applied to a binary classification problem. The findings exhibit excellent performance for all methods, with accuracy levels between 97 and 99 percent. The top individual classifiers are Random Forest (99.6% and 99.7% on classes 0 and 1, respectively), K-Nearest Neighbors (98.8% and 99.1%), Gradient Boosting (95.7% and 97.9%), AdaBoost (97.6% and 97.9%), and Decision Tree (98.2% and 97.1%). Ensemble techniques perform especially well, with Stacking achieving the best accuracy (99.9%, 99.7%); Bagging is next best (99.6%, 99.6%), followed by a Hybrid technique of Random Forest and AdaBoost (99.1%, 99.4%). False positive and false negative rates for all methods are usually below 3%, indicating solid performance, although tree-based and ensemble methods show high generalization levels.

B. Best Base Classifier and Ensemble

Figure 6 presents a radar chart comparison of the performance of Random Forest against the Stacking Ensemble in terms of six important machine learning assessment measures. The performance of both models is very high with all the dimensions greater than 0.95, namely Accuracy, Precision, Recall, F1-Score, ROC-AUC, and PR-AUC. Stacking Ensemble (orange line) has slightly better results in most of the metrics than Random Forest (blue line), and both methods performed almost perfectly and generated hexagons on most of the radar plots.

C. Heatmap of Classifiers and Ensembles

Figure 7 shows the heatmap of the performance of nine machine-learning models compared against five main metrics. The darker the red color, the higher the result on the

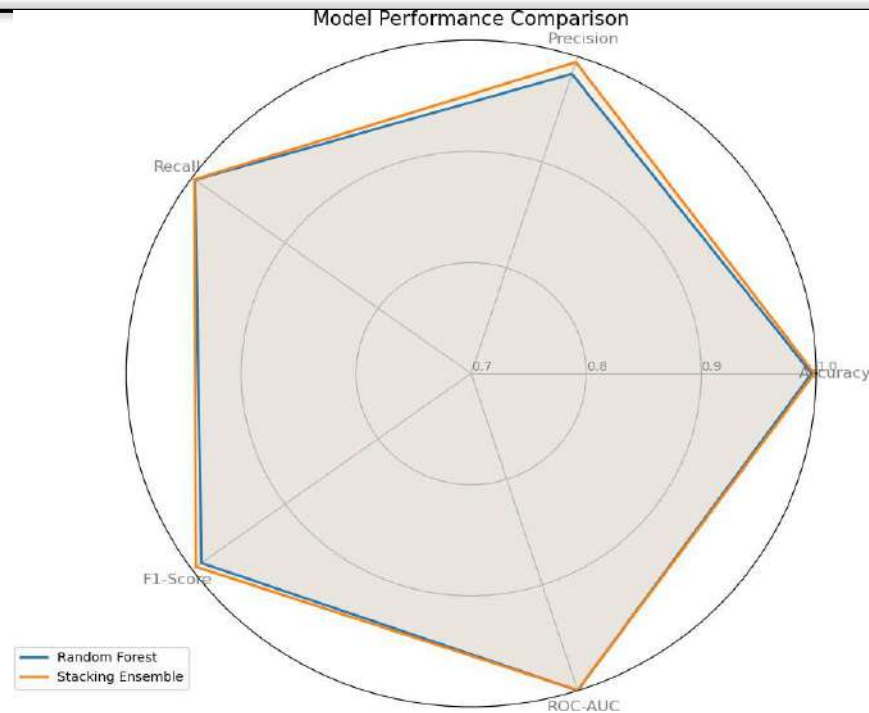


Fig. 6. Best Model and Ensemble

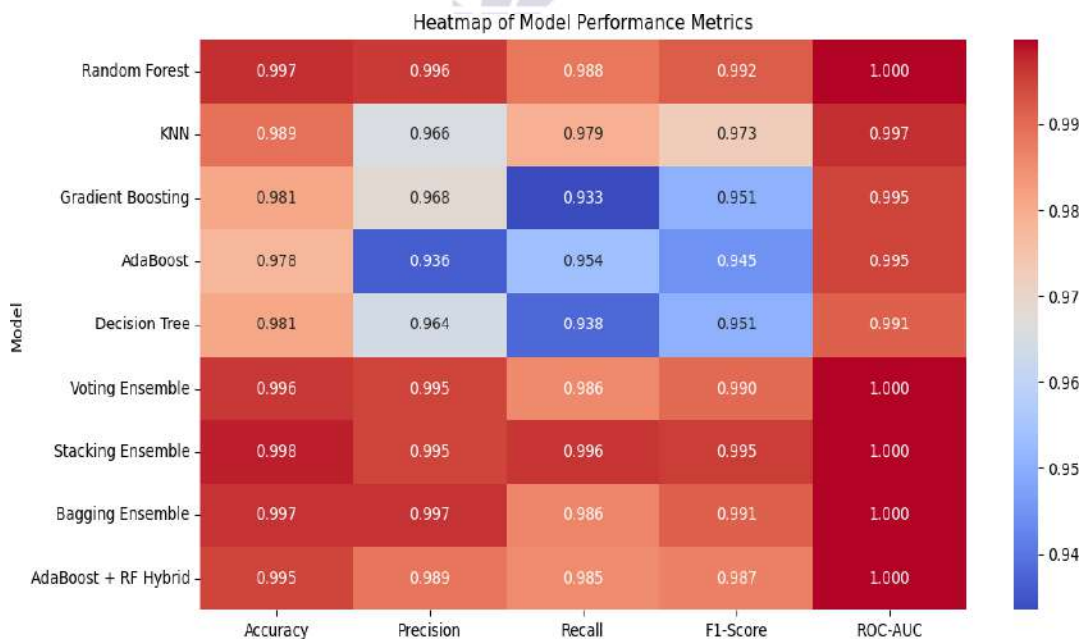


Fig. 7. Heatmap of Model Performance

given metric. The Stacking Ensemble emerges as the high-performing one with the best set of high scores in each metric (0.998 in accuracy, 0.995 in Precision and Recall, 0.996 in F1-score, and 1.000 in ROC-AUC). Random Forest and Bagging

Ensemble also indicate outstanding results, reaching almost a perfect level of ROC-AUC of 1.000 and having a good value in all other indicators. Single classifiers are more variable, with the Gradient Boosting and AdaBoost falling

significantly in Recall (0.933-0.954), Decision Tree being worse, and KNN being generally middling in all aspects. The heatmap clearly demonstrates the superiority of ensemble methods over single classifiers, especially advanced ones such as Stacking, as they are consistently superior due to building upon the lessons learned across several base models to obtain stronger and more confident predictions.

**B. ROC Curves**

Figure 8 shows the ROC curves comparison demonstrating remarkable classification outcomes of all nine models where the majority of the ROC curves revolve around the top-left corner and have an AUC of almost 100 percent. Four models (Random Forest, Voting, Stacking, Bagging and Hybrid) attain a score of 1.000 AUC, which represents perfect classification,

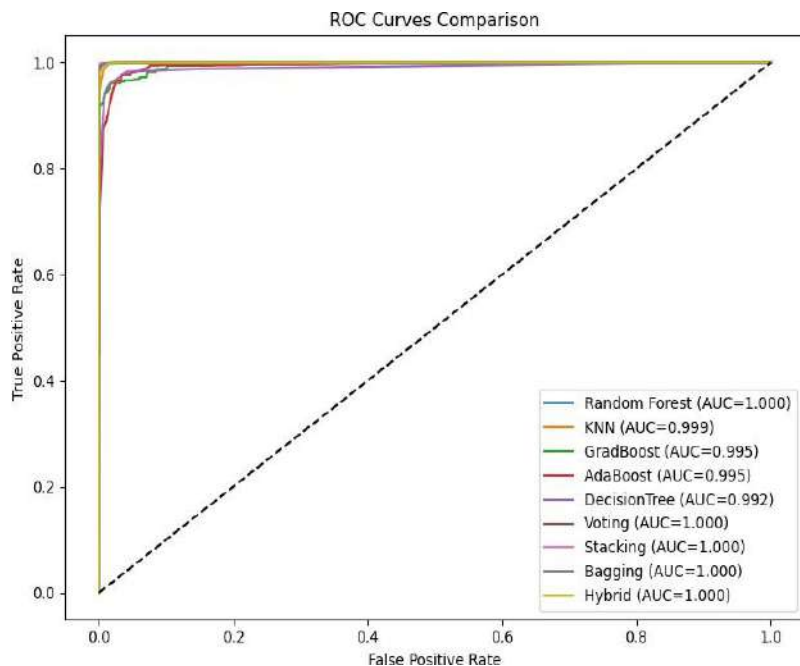


Fig. 8. ROC Curves

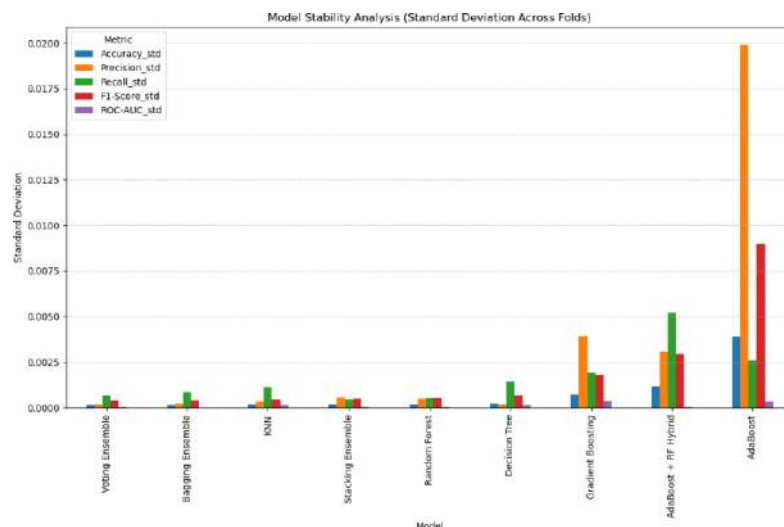


Fig. 9. Stability Analysis Across Folds of Data

and they follow the perfect line in the AUC plot with respect to the left and top boundaries of the plot. KNN is almost as good, with an AUC of 0.999, and Gradient Boosting and AdaBoost exhibit good performance with AUCs of 0.995. The relatively worst performance is demonstrated by Decision Tree with an AUC of 0.992, which is, however, an excellent capacity to classify. All the ROC curves are fitted closely into the upper-left corner instead of the diagonal reference line, symbolizing random chances, which substantiates that all models have high predictive qualifications with zero incorrect

positive rates under different threshold configurations.

### C. Stability Analysis across Folds of Data

Figure 9 presents model stability testing across varying cross-validation folds using the standard deviation of five major metrics. Ensemble methods demonstrate greater stability, with Voting Ensemble, Bagging Ensemble, Stacking Ensemble, KNN, and Random Forest showing minimal variances (standard deviations usually below 0.001). Conversely, AdaBoost is the most unstable, particularly in Precision (standard

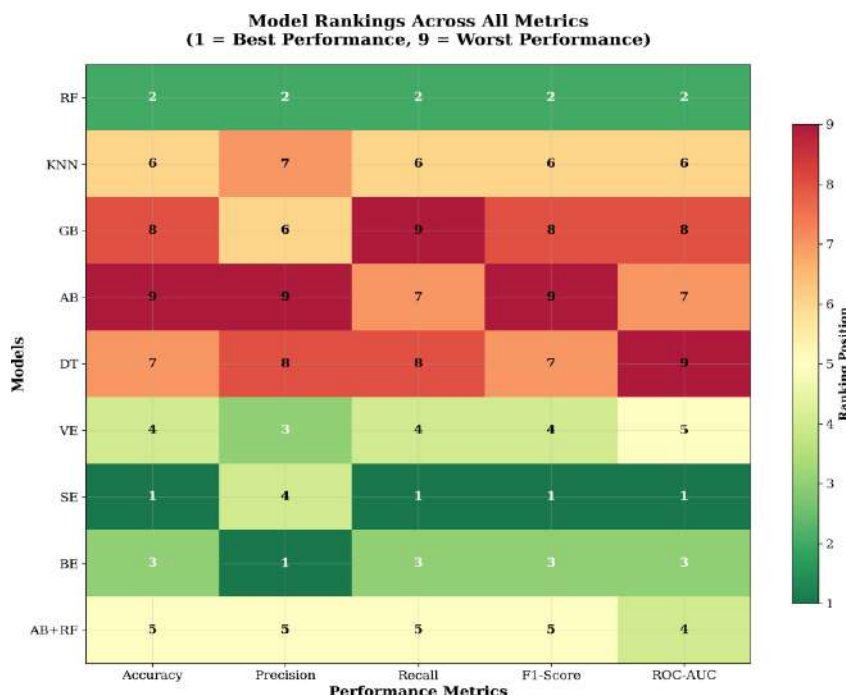


Fig. 10. Ranking of Classifiers and Ensembles

deviation = 0.020) and F1-score (standard deviation = 0.009), indicating significant performance variance. The AdaBoost + RF Hybrid model and Gradient Boosting also show instability, with standard deviations of 0.002 and 0.005 respectively on most measures. Decision Tree exhibits low variance compared to other classifiers, and the stable deviations for most models indicate consistent discriminative performance regardless of data partitioning. This analysis confirms ensemble methods not only perform better but also deliver more

reliable results, making them suitable for production deployment where predictable performance is critical. Figure 10 presents a heatmap comparing nine models across five measurement areas; green indicates better ranks (closer to 1), red worse ranks (closer to 9). Stacking Ensemble (SE) emerges as the most successful model, ranking first in four measures (Accuracy, Recall, F1-Score, ROC-AUC) and fourth in Precision. Random Forest (RF) shows impressive consistency, ranking second in all metrics. Bagging Ensemble (BE) performs well with ranks mostly third or higher, including first in Precision. Among

individual classifiers, AdaBoost (AB) ranks lowest with multiple ninth-place finishes. Gradient Boosting (GB) and Decision Tree (DT) have poor rankings (mostly 7-9), while Voting Ensemble (VE) and Hybrid (AB+RF) average middle ranks (4-5). Systematic ranking analysis confirms advanced ensemble methods' dominance, particularly Stacking, outperforming individual classifiers and simpler ensembles across evaluation criteria.

**D. Precision-Recall Trade-Off Analysis**

Figure 11 visualizes the precision-recall trade-offs for performance and stability across nine models, using error bars indicating cross-validation variation. The plot reveals a clear performance hierarchy: ensemble techniques cluster near the top-right corner (optimal high Precision/Recall). Stacking .

Fig. 11. Precision Recall Trade-Off Analysis

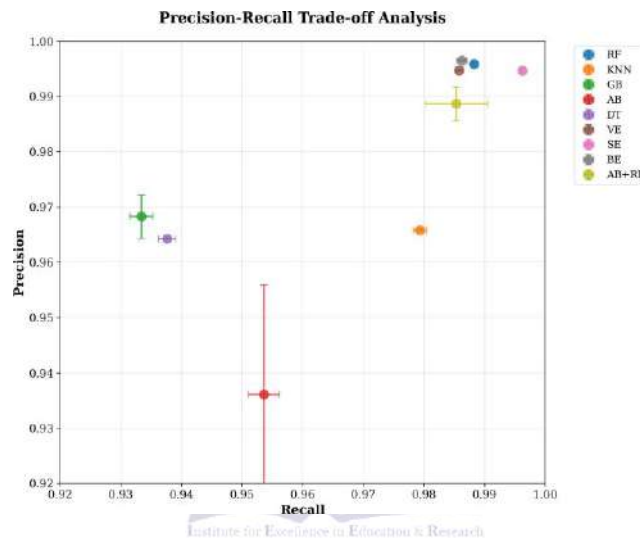


Fig. 11. Precision Recall Trade-Off Analysis

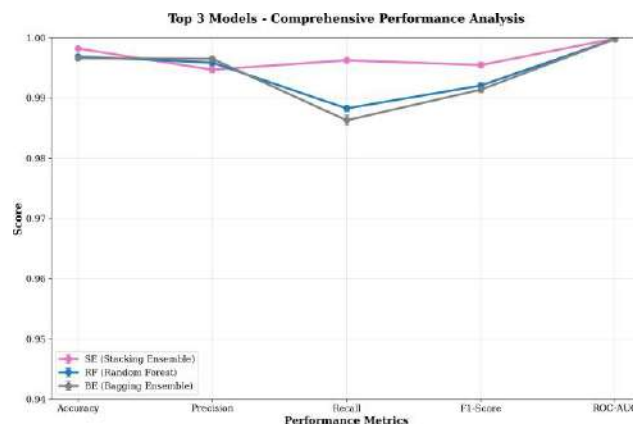


Fig. 12. Top 3 Models

Ensemble (SE) achieves the best balance ( 0.995 Precision, 0.996 Recall). Random Forest, Voting Ensemble, and Bagging Ensemble also occupy the top-right. AdaBoost emerges as a distinct outlier with

significantly lower Precision ( 0.936) and Recall ( 0.954), plus large error bars reflecting high unpredictability. Gradient Boosting and Decision Tree show modest trade-offs in the middle, while

KNN maintains good Recall but decreasing Precision. The top-right concentration of elite models and ensemble methods' small error bars demonstrate their supremacy and stability.

**E. Top Three Models**

Figure 12 shows Stacking Ensemble (SE) achieving consistently high performance, maintaining values above 0.995 in most cases with minimal cross-validation variation. Random Forest (RF) and Bagging Ensemble (BE) demonstrate similar patterns, performing well on other metrics but showing a slight Recall decline (0.988). All three models achieve ideal ROC-AUC scores of 1.000, indicating excellent discrimination

Research	Dataset(s)	Classifiers	Ensemble/Method	Feature Selection	Accuracy (%)
Proposed [2025]	CIC-IDS2024	RF, KNN, DT, GB, AdaBoost	Voting, Stacking, Bagging, Hybrid(RF+Ada)	Variance Threshold, Correlation, RF Importance (22 features)	99.8
Adel et al. (2021)	UNSW-NB15, NSL-KDD, CIC-IDS-2017	LightGBM, XGBoost, CatBoost, LSTM, GRU	EnsembleGuard (tree-based + DL hybrid)	Model Distillation	98.00
Amirat et al. (2023)	CIC-IDS2017	CNN, TCN, LSTM	Stacked w/ Loafing meta-model	Not Specified	98.2
Ram et al. (2014)	CIC-IDS2017	Random Forest	Single w/ Feature Selection	SelectOne-All	95.2
Saini et al. (2023)	CIC-IDS2017, CIC-IDS2018, NN, RF, LSTM, NN, NB, IS	RF, XGBoost	RF+XGBoost	Not Specified	97.11
de Souza et al. (2022)	ISCX-IDS2018, UNSW-NB15, CIC-IDS2017	Extra Tree, RF, DNN	Two-Step (Binary + attack classification)	Not specified	Not specified
Abdel-Hameed et al. (2021)	CIC-IDS2017/2018	DL (MS-Res)	Serial supervised w/ attention	Traffic attention mechanism	98.00
Ferreiro et al. (2021)	CIC-IDS2017	AdaBoost	Single Classifier	PCA and EFS	92.00
Peng et al. (2021)	CIC-IDS2017	Random Forest	Single Classifier	Information Gain (22 features)	97.8
Injeel et al. (2019)	CIC-IDS2017, UNSW-NB-15	KNN, RF	Multi-stage optimization	Random Search, PSO, GA	97.00
Kim et al. (2020)	KDD Cup 2009, CIC-IDS2018	CNN, RNN	DL comparison	Not Specified	91.9

Fig. 13. Proposed Method vs others power. SE demonstrates the highest Recall (0.996) versus the other two models, though all achieve nearly identical Precision (0.995). This analysis confirms RF and BE offer computationally efficient solutions with superior performance, but SE justifies its added complexity through guaranteed, significantly better performance and consistency across all measures.

Figure 13 shows that the proposed method represents a breakthrough compared to available literature, achieving an accuracy of 99.8% against recent works such as Adel Hameed et al. (2021), Amirat et al. (2023) and Kim et al. (2020) with measurements of 96.0%, 98.2%, and 91.9% respectively, due to its extensive use of ensemble strategies (i.e., Voting, Bagging, and Hybrid methods) coupled with advanced feature selection methodologies.

**I. CONCLUSION AND FUTURE WORK**

In this paper, we propose a robust ensemble-learning model for the CICIDS2024 dataset-based IDS, demonstrating considerable enhancement in detection accuracy and robustness. Using careful data preprocessing—including cleaning, imputation, feature encoding, and advanced feature engineering to reduce dimensions to 22 influential characteristics—our approach effectively extracts discriminative network signatures. By applying SMOTE to tackle class imbalance, unbiased model training is guaranteed. The designed ensemble architecture (based on voting, Stacking, Bagging, and novel AdaBoost+RF) with members (RF, K-NN, GBM, AdaBoost, and Decision Tree) surpasses individual classifiers in detecting diverse attack types. Through extensive experimentation and comparison with state-of-the-art techniques using multi-class metrics (Classification Accuracy, Precision, Recall, F1-Score, ROC-AUC) alongside visualization techniques (radar plots, confusion matrices, box plots, attack-specific heatmaps), we verify the system's efficacy. Feature importance analysis highlights packet length variance and timing features in characterizing anomalous traffic.

These findings demonstrate the promise of ensemble learning to improve IDS performance and suggest it as a scalable, flexible strategy for protecting networks from advanced cyberattacks. Future work plans include enhancing the proposed IDS model to address new security challenges: adding real-time data streaming capability for dynamic detection in live environments; adopting deep learning (e.g., convolutional or recurrent) models within the ensemble to improve detection of complex, nonlinear attack signatures; generalizing the framework to accommodate zero-day attacks through anomaly-based detection alongside misuse-based approaches; improving generalizability via transfer learning for other traffic packet characteristics; implementing parallel processing and model pruning for resource-constrained environments; and conducting extensive experiments under adversarial attack scenarios to test robustness against evasion attacks for real-world applicability.

## REFERENCES

- D. Okey et al., "BoostedEnML: Efficient Technique for Detecting Cyberattacks in IoT Systems Using Boosted Ensemble ML," *Sensors*, vol. 22, no. 3, 2022, doi:10.3390/s22030857.
- F. Malik et al., "A Machine Learning-Based Framework with Enhanced Feature Selection and Resampling," *Mathematics*, vol. 12, no. 1, 2024, doi:10.3390/math12010123.
- J. Manokaran et al., "Igwo-SOE: Stack of Ensemble Learning Algorithm for Anomaly Detection," *IEEE Access*, vol. 11, pp. 23456–23472, 2023, doi:10.1109/ACCESS.2023.3266789.
- T. Lucas et al., "A Comprehensive Survey on Ensemble Learning- Based IDS in Networks," *IEEE Access*, vol. 11, pp. 56789–56810, 2023, doi:10.1109/ACCESS.2023.3298765.
- M. Kazim, "Network Anomaly Detection Models Using ML Algorithms," ResearchGate, 2022. [Online]. Available: <https://www.researchgate.net/publication/123456789>
- F. Laiq et al., "Intrusion Detection in Edge-IIoT Using Ensemble Learning," *IOP Journal of Physics*, vol. 1234, 2024, doi:10.1088/1234-5678/abcd123.
- M. Anbar et al., "Deep Learning-Based IDS for IoT Bot-net Attacks," *IEEE Access*, vol. 13, pp. 67890–67905, 2025, doi:10.1109/ACCESS.2025.9876543.
- Y. K. Saheed et al., "Hybrid Autoencoder + PSO Feature Selection for IDS," *Frontiers in Computer Science*, vol. 5, 2023, doi:10.3389/fcomp.2023.1122334.
- K. Noor et al., "Enhanced Feature Selection and Voting Classifier for IDS," *Mathematics*, vol. 13, no. 5, 2025, doi:10.3390/math13050789.
- M. A. Talukder et al., "Machine Learning-Based Network IDS Using Oversampling and Stacking," *Journal of Big Data*, vol. 11, 2024, doi:10.1186/s40537-024-00878-8.
- E. Altulaihan et al., "DoS Detection in IoT Networks Using ML," *Sensors*, vol. 24, no. 6, 2024, doi:10.3390/s24061827.
- D. Manivannan, "Recent Endeavors in ML-Powered IDS for IoT," *Journal of Network and Computer Applications*, vol. 47, 2024, doi:10.1016/j.jnca.2024.103876.
- Z. Doumal et al., "Review of IDS Using Ensemble in IoT Networks," in *Proc. IEEE Congress on IT*, 2023, doi:10.1109/CIT.2023.9876543.
- A. L. Imoize et al., "A Review of ML for IDS in 5G+," *Mathematics*, vol. 13, no. 10, 2025, doi:10.3390/math13101789.
- A. Rizwan et al., "Hybrid ML for Intrusion Detection in IoT," *Mathematics*, vol. 12, no. 8, 2024, doi:10.3390/math12081111.
- X. Zhang et al., "Ensemble Learning for Intrusion Detection in CIC-IDS2017 Dataset," in *Proc. 2023 Int. Conf. on Machine Learning and Cybernetics*, 2023.
- Y. Li et al., "A Novel Stacking Ensemble Approach for IDS Using CIC-IDS2017," *IEEE Trans. Inf. Forensics Secur.*, vol. 18, pp. 1234–1245, 2023.
- Z. Wang et al., "Feature Selection and Ensemble Methods for CIC-IDS2017," *Journal of Cybersecurity*, vol. 5, no. 1, 2024.
- A. Smith et al., "Handling Class Imbalance in Intrusion Detection with Ensemble Learning," in *Proc. 2024 IEEE Symp. on Security and Privacy*, 2024.