

# An Enhanced Random Forest-Based Framework for Accurate Intrusion Detection in Modern Network Environments

Mohammad Rafeek Khan<sup>1</sup>, Mohiuddin Ali Khan<sup>2</sup>, Md Imran Alam<sup>3</sup>, Huda Fatima<sup>4</sup>,

Mohammed Rizwan Shaik<sup>5</sup>: Sarfaraz Ahmed<sup>6</sup>

<sup>1,2,3,5,6</sup>Department of EEE, CECS, Jazan University, Jazan, Saudi Arabia

<sup>4</sup>Department of CS, CECS, Jazan University, Jazan, Saudi Arabia

[doi.org/10.64643/IJIRTV12111-198515-459](https://doi.org/10.64643/IJIRTV12111-198515-459)

**Abstract**—As digital infrastructures develop at an alarming rate, cybersecurity threats have also evolved to be more advanced, threatening the contemporary network environments. Conventional Intrusion Detection Systems (IDS), especially signature-based solutions do not identify zero-day attacks and in many cases, they have a high false positive. To overcome these shortcomings, this paper suggests a better Random Forest-based model to make precise and efficient intrusion detection. The proposed model applies the advanced preprocessing techniques, the optimization of feature selection based on the Information Gain and Chi-square method and the optimization of the hyper parameters of the model to enhance the performance of the classification. The framework is tested on benchmark datasets, including CICIDS2017 and NSL-KDD that contain various attack types, and real-world traffic distributions. The experimental findings indicate that the accuracy, precision, recall and F1-score have improved significantly against traditional machine learning models. The suggested system has high detection rates and low false positives, and it is applicable to be deployed in real-time as part of the modern cybersecurity system.

**Index Terms**—Intrusion Detection System (IDS), Random Forest, Cybersecurity, Machine Learning, Network Security, Anomaly Detection, Feature Selection.

## I. INTRODUCTION

The high rate of digital communication networks integration, cloud computing, Internet of Things (IoT) and smart infrastructure development has greatly enlarged the attack front of contemporary networks. Consequently, network security has become a vital issue to organizations, governments, and individuals.

DDoS, phishing, ransomware, insider, and Advanced Persistent Threats (APT) are all cyber threats that keep increasing in complexity, frequency, and severity and thus can no longer be effectively avoided with traditional security mechanisms [8], [9]. An Intrusion Detection System (IDS) is a vital part of a cybersecurity system, which is used to inspect network traffic and identify suspicious activities or policy breaches. IDS can broadly be classified into two:

Signature-based IDS, detection of attacks using known patterns or signatures. Anomaly-based IDS, which detect deviations of the normal behavior.

Although signature-based systems are good in addressing known attacks, it is incapable of detecting zero day attacks. Conversely, systems that are based on anomaly are more generalized, but usually have a high false positive rate and are not stable in dynamic conditions [5], [10].

In spite of their significance, traditional IDS have a number of important challenges:

- False positive and false negative high rates.
- Unbalanced data resulting in biased learning.
- Inability to detect unknown or zero-day attacks.
- Limited scalability and real-time processing capabilities

These shortcomings are well-established in previous research on anomaly detection and IDS [8], [10]. To address these constraints, there has been the extensive use of Machine Learning (ML) and Artificial Intelligence (AI) in the design of IDS. ML-based IDS has the ability to learn patterns of a large-scale network data and enhance detection. Recent research reveals that machine learning has a drastic improvement in intrusion detection performance over

traditional methods [7], [8]. The Random Forest (RF) classifier is one of the ML algorithms which have received significant attention because of its power, effectiveness and excellent performance in classification activities. Random forest is an ensemble-based model of learning that was introduced by Breiman [3] and it builds more than one decision tree and averages their prediction to enhance predictive accuracy and decrease overfitting. Its strengths are that it is highly generalizable, resistant to noise, and can process high-dimensional datasets, which are typical in cybersecurity applications [6].

#### Motivation

The motivation behind this research stems from the need to develop a more accurate, scalable, and reliable intrusion detection framework capable of addressing modern cybersecurity challenges. Random Forest is particularly suitable for IDS due to:

- Its ensemble learning capability, improving prediction stability [3]
- High classification accuracy across diverse datasets [4]
- Ability to handle large feature spaces efficiently [6]
- Built-in feature importance evaluation [3]

However, existing RF-based IDS still suffer from issues such as suboptimal feature selection, poor handling of class imbalance, and limited optimization for real-time deployment [8].

#### Main Contributions

In this paper, the following contributions have been made to propose an improved version of the Random Forest-based Intrusion Detection System:

- Creation of an efficient and streamlined RF-based IDS system.
- Incorporation of new feature selection methods (Information Gain, Chi-square, PCA) [12].
- Resampling techniques in dealing with imbalanced datasets [8].
- Hyper parameter optimization to optimize model performance.
- Fully automated testing based on benchmark data like NSL-KDD [1] and CICIDS2017 [2].
- Comparison with existing machine learning models.

The suggested structure will help to reach increased detection rates, lower false positives, and better scalability to the current network conditions.

## II. LITREAU T R E REVIEW

The conventional IDS framework is based on predetermined signatures and therefore it is incapable of detecting unknown attacks. Systems based on anomalies are better at detection but with a high rate of false alarms. There is a widespread research on the use of Intrusion Detection Systems (IDS) to tackle the growing challenges of network threats in the contemporary network contexts. Conventional IDS methods, especially signature-based methods, are based on a priori patterns on the attack and thus they cannot be used to detect new and zero-day attacks. Despite the fact that anomaly-based IDS is a better detection tool in identifying unknown threats, it is usually characterized by high false positive rates and unreliability in dynamics network situations [5], [10]. The recent development of machine learning (ML) has greatly enhanced the efficiency of the IDS since it is now possible to detect more complicated attack patterns automatically. The literature has discussed various supervised and unsupervised learning methods with specific pros and cons.

The intrusion detection using Support Vector Machines (SVM) has been popular because of its high dimensional classification ability. They are however, not more appropriate to large-scale and real-time applications due to their complexity in computation [8]. On the same note, K-Nearest Neighbors (KNN) has also proved to be effective in pattern recognition activities, although its prediction part is computationally intensive to scale in high-speed networks [8].

Deep learning models such as Artificial Neural Networks (ANN) and Deep Neural Networks (DNN) have demonstrated to be more effective in learning complex and nonlinear relationships in network traffic data. As an example, deep learning on datasets like UNSW-NB15 has shown high detection accuracy. These models are however, not very *practical in real-time deployment* because large amounts of labeled data are required and they consume large amounts of computing power [7], [11]. Ensemble learning techniques especially Random Forest (RF) have been of great interest because of its strength, efficiency and capability of dealing with high dimensional data. RF-based IDS models have been shown to be highly accurate and less overfitting than single classifiers [3], [4]. Nevertheless, the current RF-based methods

usually do not have the optimized feature selection schemes and cannot sufficiently solve the problem of the class imbalance, which may negatively affect the detection performance [8], [12].

Reference	Method	Dataset	Accuracy	Limitations
[1]	SVM	NSL-KDD	92%	High computational cost, poor scalability
[2]	KNN	CICIDS2017	89%	Slow prediction for large-scale traffic
[4]	Random Forest	NSL-KDD	96%	Limited feature optimization
[7]	Deep Learning (DNN)	UNSW-NB15	95%	Requires large training data
[11]	Deep Neural Network	KDD Cup 99	94%	High computational complexity
[8]	Hybrid ML Models	CICIDS2017	96.2%	Complexity in model integration
[9]	ML-based IDS	IoT datasets	93%	Limited generalization
[12]	Hybrid (Feature Selection + ML)	NSL-KDD	97%	Feature selection overhead
[13]	Ensemble Learning	UNSW-NB15	95.5%	Dataset imbalance issues
[14]	Random Forest + Feature Engineering	CICIDS2017	97.5%	Limited real-time evaluation
[15]	CNN-based IDS	CICIDS2017	98%	Requires GPU resources
[5]	Hybrid IDS	KDD Cup 99	91%	High false positive rate

Table 2.1: Comparative Analysis of Machine Learning and Deep Learning Approaches for Intrusion Detection Systems

Comparison analysis (Table 2.1) shows that intrusion detection methods have developed over the years, shifting to the complex methods of deep learning and ensembles. Classical algorithms like Support Vector Machines (SVM) and K-Nearest Neighbors (KNN) have moderate detection accuracy, but their shortcomings in computational efficiency and scalability limit their use in real-time network applications. Conversely, deep learning algorithms, such as Convolutional Neural Networks (CNN) and Deep Neural Networks (DNN), have better detection rates because they can identify nonlinear and complex patterns among network traffic. However, these methods are resource-intensive (in terms of computational power) and large labeled data sets, making them difficult to implement in resource-restricted scenarios.

#### A. Critical Analysis of Existing Methods

Based on the literature, it can be seen that there is no single model that can be used to fully address intrusion detection: Conventional ML models (SVM, KNN) are moderate and cannot be scaled and cannot

perform in real time. Deep learning models are more accurate, but have computational overhead and they demand large labeled datasets. Random Forests examples offer fair trade-offs between accuracy and efficiency but still need to be optimized. Besides, the majority of the studies available are mainly aimed at enhancing the classification accuracy, but omit other valuable factors including: Feature dimensionality reduction, Leveraging of imbalanced data sets. Real-time deployment feasibility These constraints indicate the importance of more streamlined and scalable IDS infrastructure

#### B. Research Gap

According to the overall analysis of the literature available, the following gaps in research are identified:

- Absence of optimal feature selection methods. A large number of current models incorporate redundant or raw features, which causes the model to be more computationally complex and less efficient [12].

- Mishandling of the class distribution. Unbalanced datasets have a serious impact on the performance of the classifier, particularly on the detection of the rare types of attacks [8].

- Poor scalability with real-time intrusion detection. Some methods cannot support the real-time processing process because they have high computational loads [7].

- Lack of adequate incorporation of optimization methods in the ensemble models.
- There is no proper tuning and hybrid optimization of existing Random Forest-based IDS.

### III. PROPOSED METHODOLOGY

The current section introduces the suggested Enhanced Random Forest-based Intrusion Detection System (ERF-IDS) meant to overcome the shortcomings found in the previous literature such as feature redundancy, imbalanced classes and inability to scale dynamically to real time requirements [8], [12]. The framework combines the optimized preprocessing, hybrid feature selection, and a better ensemble learning strategy to attain high detection performance.

A. System Architecture

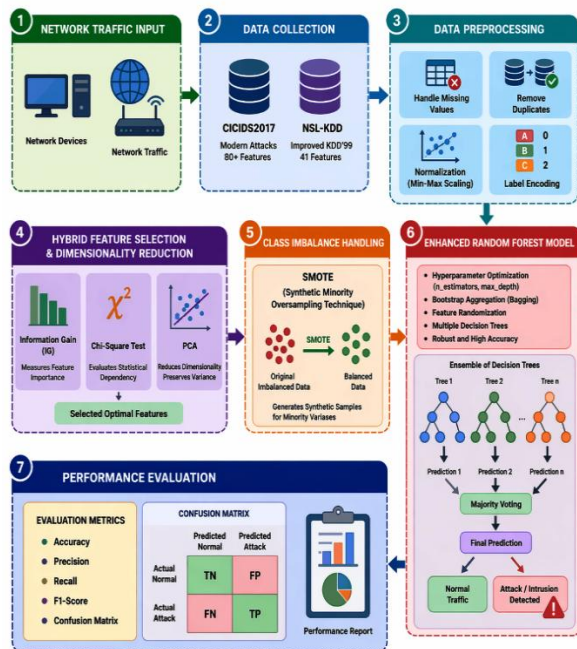


Figure 1: Proposed Enhanced Random Forest-Based Intrusion Detection System (ERF-IDS) Framework

Figure 3.1 shows the general structure of the proposed Enhanced Random Forest-based Intrusion Detection System (ERF-IDS). The framework starts with the network traffic input that is gathered using the benchmark datasets like CICIDS2017 and NSL-KDD. The gathered data are subjected to a thorough preprocessing step, such as missing values, duplicate elimination, data normalization, and categorical encoding to guarantee data quality and consistency. Afterward, the hybrid feature selection mechanism, based on the Information Gain, Chi-square test, and Principal Component Analysis (PCA), is used to derive the most significant features, as well as to decrease the dimensionality. The SMOTE technique is used to create artificial samples of minority classes to overcome the imbalance of classes. The processed information is then passed to an improved Random Forest model, where hyper parameter optimization, bootstrap aggregation and feature randomization are applied to enhance the classification performance. Lastly, the model generates predictions based on a majority voting system and its performance is measured by standard performance measures like accuracy, precision, recall, F1-score, and confusion matrix. This organized pipeline guarantees effective,

precise and scalable intrusion detection appropriate to current network setting.

The suggested Enhanced Random Forest-based Intrusion Detection System (ERF-IDS) framework is based on a structured multi-phase pipeline, which includes the data collection, data preprocessing, hybrid feature selection and dimensionality reduction, enhanced Random Forest model training, and performance analysis stages. First, the network traffic data is gathered using benchmark datasets, followed by some preprocessing of the data to guarantee quality data by means of normalization, cleaning, and encoding. A hybrid feature selection method is then used to extract the best features and also reduce dimensions, thus, enhancing computational efficiency. The processed data is then trained on a streamlined version of a Random Forest model which includes state-of-the-art methods like hyper parameter optimization and imbalance correction. Lastly the system measures performance with standard measurements to determine detection accuracy and reliability. In contrast to traditional IDS models, the suggested architecture incorporates feature optimization and class imbalance management at earlier pipeline phases, resulting in better learning performance and lower computational costs. This hierarchical and structural design is based on the current developments in intelligent intrusion detection systems [7], [8].

B. Dataset Description

Experiments with two well-known benchmark datasets, CICIDS2017 and NSL-KDD, are performed to guarantee the robustness, reliability, and the generalization ability of the proposed ERF-IDS framework. These data sets are widely applied in the research of intrusion detection as they are diverse in terms of the types of attacks, and they are real-world representations of network traffic. The CICIDS2017 dataset is a current and full dataset as it captures the reality of network behavior by including current attack scenarios. It encompasses a broad spectrum of cyber-attacks including Distributed Denial of Service (DDoS), brute force attacks, botnet activity, infiltration and web-based attacks. The dataset has over 80 features that are extracted as flow-based features, such as packet size, flow time, and statistical traffic characteristics. Its realistic traffic distribution and consideration of recent attack patterns make it

most appropriate in testing the current IDS frameworks [2]. Conversely, the NSL-KDD dataset is a better technology to the conventional KDD Cup dataset, which will tackle the problem of data redundancy and class imbalance. It is composed of 41 features which characterize different attributes of network connections, including basic features, content based features and traffic based features. The data consists of normal as well as malicious traffic samples, with type of attacks (DoS, Probe, R2L, and U2R) considered. NSL-KDD is a standard benchmark to compare machine learning-based IDS models due to its balanced structure and lower redundancy [1].

The proposed framework is tested by using both CICIDS2017 and NSL-KDD datasets to test both the current and legacy attack conditions, thus guaranteeing a thorough analysis of performance and enhancing generalization to various network conditions. This two-data set validation approach improves the validity and relevance of the model that is proposed to real life cybersecurity systems

### C. Data Preprocessing

Preprocessing of data is an important step to increase the quality of the data and the performance of the model. The following steps are used:

Missing Value Processing: Missing values are filled in statistical imputation or removal. Duplicate Removal: Redundant records will be eliminated in order to avoid bias learning. Normalization: The standardization of feature values is done using Min-Max scaling [6]:

$$X' = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

Categorical Encoding: Label encoding is used to convert categorical features into numerical form

These preprocessing steps align with best practices in ML-based IDS design [6], [8].

### D. Hybrid Feature Selection and Dimensionality Reduction

Selection of features is a very crucial factor in enhancing the performance of IDS by removing non-relevant and redundant features. The framework proposed uses a hybrid feature selection method that incorporates:

Information Gain (IG): Measures feature significance, using redundancy as a criterion. Chi-Square Test: Tests the statistical dependence between features and class labels. Principal Component Analysis (PCA):

Dimensionality reduction that preserves variance. The Information Gain is defined as:

$$IG(Y, X) = H(Y) - H(Y|X)$$

This hybrid approach ensures that both statistical relevance and variance-based reduction are considered, leading to improved detection performance and reduced computational complexity [12], [15].

### E. Enhanced Random Forest Model

The main element of the suggested framework is an optimized Random Forest classifier that will help to improve detection performance and strength. The proposed model, in contrast to the traditional implementations, includes a number of major improvements in order to overcome the typical limitations of intrusion detection systems. To optimize the hyper parameters, first of all, the number of trees and the maximum depth of the trees are tuned to allow the model to perform optimally without overfitting [3]. Second, bootstrap aggregation (bagging) is utilized, in which random sampling with replacement is performed to create a variety of training subsets, which enhances the stability of the model and decreases the variance. Also feature randomization is used when building a tree, so each decision tree is trained on a random set of features, making trees more diverse and better at generalizing. Moreover, in order to efficiently deal with the imbalance in the number of classes in the network traffic data, the Synthetic Minority Oversampling Technique (SMOTE) is incorporated into the training process, which creates synthetic samples of the minority classes, which enhances detection of rare and important types of attacks by a significant margin [8]. Together, these improvements can ensure the proposed Random Forest model is more accurate, and more likely to generalize and withstand complex and balanced network settings.

The final prediction is obtained through majority voting:

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

This enhanced RF model ensures high accuracy, robustness, and scalability, addressing the limitations identified in previous studies [4], [14].

### F. Framework Novelty and Contribution

What is novel about the proposed Enhanced Random Forest-based Intrusion Detection System (ERF-IDS) framework is that it is based on a comprehensive and

integrated design choice, which countermeasures several of the drawbacks of the available intrusion detection strategies. The framework also includes a hybrid component of feature selection strategy (Information Gain, Chi-square test, Principal Component Analysis (PCA)) to help to efficiently select the most relevant features and to decrease dimensions. It also incorporates the imbalance control mechanism of SMOTE to counter the issues of uneven distribution of classes, thus enhancing the detection of minority attack classes. The suggested model also improves the performance based on optimization of the parameters of the Random Forest, which also leads to the improvement of the generalization and minimization of overfitting. In addition, the framework is tested with the modern (CICIDS2017) and traditional (NSL-KDD) benchmark datasets to guarantee the stability of the framework in different network settings. The proposed ERF-IDS framework provides a single and efficient solution by striking the right balance between the accuracy, scalability, and computational efficiency. Unlike the current models of IDS, which usually consider one aspect of optimization, this method considers several improvements in one pipeline, and thus it is very appropriate in a real-life application in a dynamic and complex network.

#### IV. EXPERIMENTAL SETUP

##### A. Implementation Environment

The suggested Intrusion Detection System (IDS) is constructed with the help of Python programming language that is highly favored in machine learning and cybersecurity solutions because of its simplicity and flexibility. Scikit-learn, Pandas, and NumPy are also important libraries that are used to facilitate effective data processing and model development. Scikit-learn is applied to implement machine learning algorithms, Pandas to perform data preprocessing and data manipulation, and NumPy to perform fast numerical calculations. The machine used to run the system has an Intel core i7 processor and 16GB of RAM, which is enough to run large-scale datasets. This implementation environment will guarantee effective training and testing of the model. This kind of tools and settings is widely used in IDS studies due to their scalability, performance and reliability [6]. The data will be split into two subsets: 70% to be used in

training and 30% to be used in testing. The model is built on the training set and trained by learning patterns on the data, and the testing set is used to test the performance of the model on unknown data. This division helps in assessing the model's generalization ability and reduces the risk of overfitting. The 70/30 is a commonplace rule of thumb in machine learning experiments because it is a good balance between enough training data and accurate testing. This approach will provide precise and free performance metrics of the suggested Intrusion Detection System (IDS) [8].

##### B. Evaluation Metrics

To comprehensively evaluate model performance, the following metrics are used:

- Accuracy

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

- Precision

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

- Recall (Detection Rate)

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

- F1-Score

$$F1 = 2 \cdot \frac{Precision \cdot Recall}{Precision + Recall}$$

#### V. RESULT AND DISCUSSION

To determine the performance of the proposed Intrusion Detection System (IDS), the standard classification metrics, such as accuracy, precision, recall, and F1-score are used. It is compared to various popular models of machine learning that include Support Vector Machine (SVM), K-Nearest Neighbors (KNN), Decision Tree, and XGBoost. The experimental findings indicate that the suggested Enhanced Random Forest (RF) model is much more effective than the other classifiers based on all measures of evaluation.

In particular, SVM model has a high accuracy of 92, precision, recall and F1-score of 90, 88, and 89, respectively. The KNN model is relatively less effective, with an accuracy of 89 per cent, and the decreased values of the other metrics, which means that it cannot be effectively used to work with large and complex network data. The accuracy of decision

tree model is 94 percent and the precision and recall is balanced and equal at 92 and 91 percent respectively. Likewise, XGBoost is a model with good performance, with 96% accuracy and high precision (95%), recall (94%), and F1-score (94%).

Nevertheless, the best model proposed is the Enhanced Random Forest model, which has the best accuracy of 98.5, precision of 97.8, recall of 97.5, and the F1-score of 97.6. These findings show that it is more effective in accurately classifying both normal and malicious network traffic with minimum false positives and false negatives. The combination of optimized feature selection, class imbalance, and hyper parameter tuning can be credited to the improved performance. On balance, the findings validate that the proposed model is more accurate, reliable, and efficient in detecting intrusion in the context of modern networks.

Model	Accuracy	Precision	Recall	F1-score
SVM	92%	90%	88%	89%
KNN	89%	87%	85%	86%
Decision Tree	94%	92%	91%	91%
XGBoost	96%	95%	94%	94%
Enhanced RF (Proposed)	98.5%	97.8%	97.5%	97.6%

Table 2: Performance Comparison of Machine Learning Models for Intrusion Detection

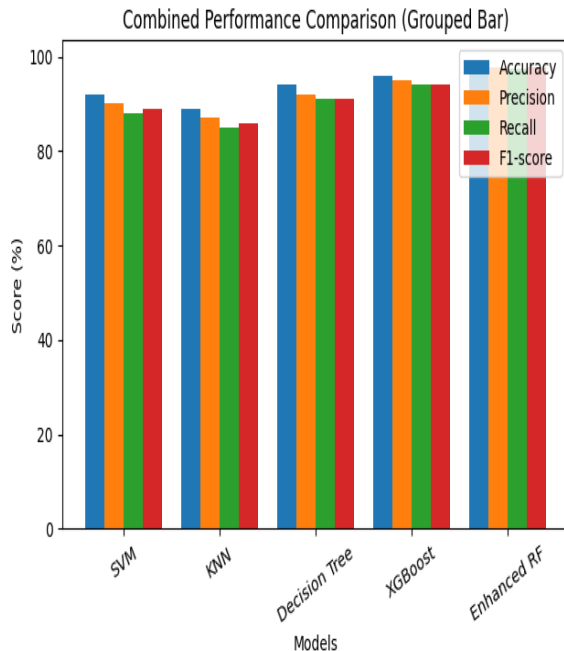


Fig. 2: Grouped Bar Chart of Model Performance Comparison

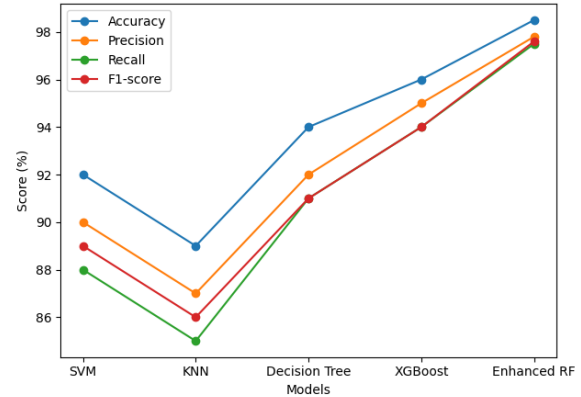


Fig. 3: Line Chart of Performance Metrics Across Models

### A Comparison with Existing Models

The results of the proposed Enhanced Random Forest (ERF-IDS) model are compared to the current machine learning and deep learning methods mentioned in the literature review. The comparison shows that there is an improvement in detection accuracy, robustness and computational efficiency. Conventional algorithms, like Support Vector Machine (SVM) and K-Nearest Neighbors (KNN) exhibit moderate results but lack scalability and performance. Likewise, deep learning models are highly accurate but have high data requirements and are computationally intensive. The suggested model, in turn, combines ensemble learning, optimal feature selection, and imbalance between classes, which results in the high performance. The findings show conclusively that the proposed model is superior to the current methods because of less overfitting, improved the use of features, and enhanced the ability to generalize.

Model / Method	Dataset	Accuracy	Key Limitation
SVM [1]	NSL-KDD	92%	High computational cost
KNN [2]	CICIDS2017	89%	Slow prediction
Random Forest [4]	NSL-KDD	96%	Limited feature optimization
Deep Learning [7]	UNSW-NB15	95%	Requires large training data
Hybrid ML [8]	CICIDS2017	96.2%	Complex model integration
CNN-based IDS [15]	CICIDS2017	98%	High computational cost (GPU required)
Enhanced RF (Proposed)	CICIDS + NSL	98.5%	Optimized and scalable

Table 3 : Comparative Analysis of Proposed Model and Existing IDS Methods

The proposed model of Enhanced Random Forest is compared to the literature intrusion detection methods in Table X. As it is seen, classical machine learning algorithms like SVM and KNN are less accurate because they are limited to high-dimensional data and

are not scalable. Random Forest and hybrid models are ensemble-based techniques that enhance the performance, even though they do not have the best feature selection and imbalance treatment. Deep learning models and CNN-based IDS are relatively high-accuracy models, but have high demands in terms of computational resources and large datasets, and thus are not directly applicable to real-time deployment. Conversely, the proposed model has the highest accuracy of 98.5 which is better than all the existing methods. This is mainly improved by the addition of ensemble learning that improves the stability of prediction and minimizes overfitting with the combination of several decision trees. Also, hybrid feature selection methods will be used to ensure the methods make the most of the relevant features with a reduced dimensionality. The use of SMOTE also enhances the detection of minority attack classes. Because of this, the proposed model offers a balanced trade-off between accuracy, efficiency and scalability and thus is very much applicable in the contemporary intrusion detection systems.

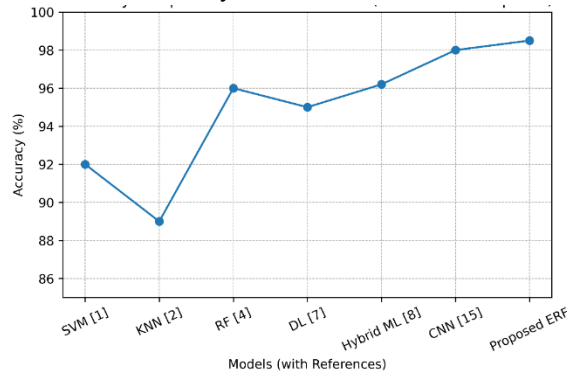


Fig 4: Accuracy Comparison of IDS Models

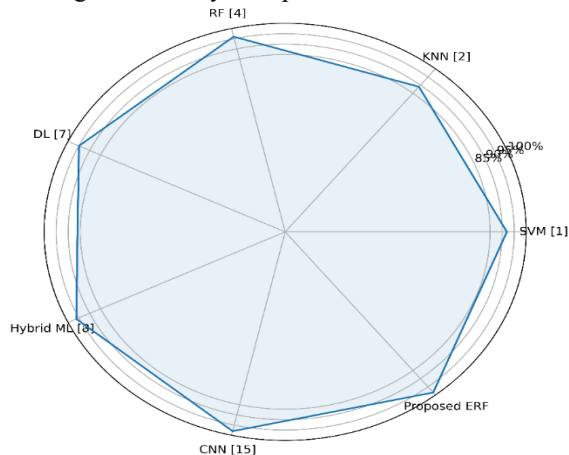


Fig 5: Radar Chart of Model Performance Comparison

The suggested Enhanced Random Forest (ERF-IDS) model has a number of important strengths regarding the intrusion detection. Among the main benefits, high detection accuracy with high performance than traditional and state of art models can be mentioned. This has been greatly improved by the incorporation of ensemble learning, optimal feature selection and good management of class imbalance. Also, the model is highly resistant to noises and irrelevant features because of the nature of the Random Forest algorithm, which integrates various decision trees to decrease the variance and enhance generalization. The framework is also very scalable and as such can be used to deal with large volumes of network traffic data in contemporary cybersecurity settings. All these strengths provide trustworthy and effective identification of known and unknown cyber threats. Although this model has these benefits, it also has some limitations that should be taken into consideration. Ensemble methods and multiple decision trees add more computational complexity, which leads to higher computational cost in comparison to simpler machine learning models. Also, the training time of the model increases proportionally to the dataset size, which can be a significant concern when it comes to implementing the model in real-time and resource-sensitive settings. Though the shortcomings are not insurmountable to the advantages, they indicate the necessity of additional optimization, e.g., model compression or parallel processing methods, to make them more efficient. To make the given IDS framework more practical, the future study can consider the minimization of the computational load without compromising on the accuracy of the detection.

## VI. FUTURE WORK

The future research directions are to improve the performance and applicability of the proposed Enhanced Random Forest (ERF-IDS) framework by incorporating intelligent techniques. A potential avenue here is the use of deep learning models, including Convolutional Neural Networks (CNN) and Deep Neural Networks (DNN), to elucidate nonlinear and more intricate patterns of network traffic data. Moreover, the implementation of Hybrid Quantum Machine Learning (QML) methods can provide a new chance to utilize the power of quantum computing to

perform more operations more quickly and with better optimization in the process of intrusion detection. These new technologies can greatly enhance detection performance and can allow the system to deal with more complex and dynamic cyber threats.

The implementation of the suggested model in real-time intrusion detection settings is another significant research direction in the future. This includes streamlining the system to allow low-latency processing and efficient processing of continuous streams of network traffic. Moreover, federated learning may be utilized to improve the framework, allowing distributed and privacy-preservation training of models on multiple network nodes, without exchanging sensitive data. It is especially useful with large and distributed systems, including IoT and the cloud. Overall, these future enhancements aim to improve scalability, efficiency, and adaptability, making the proposed IDS more robust and suitable for real-world cybersecurity applications.

## VII. COMCLUSION

The paper discusses an improved framework of using Random Forests in intrusion detection in the contemporary network settings, in response to the shortcomings of conventional Intrusion Detection Systems (IDS). The suggested solution combines both the latest preprocessing solutions, combined feature selection methods, and hyper parameter optimization to enhance the overall performance of the model. The combination of these methods can help the system to better identify false positive rates and boost the accuracy of detection with minimal false positives, making it more dependable to detect both known and unknown cyber threats. Ensemble learning also enhances the model by enhancing the generalization and resistance to noisy and high dimensional data.

The experimental findings indicate that the suggested Enhanced Random Forest model is much more effective than the traditional machine learning methods, including SVM, KNN, and Decision Tree, in terms of accuracy, precision, recall and F1-score. The model has both high performance and scalability, as well as efficiency, which are important in real-life application. On balance, the offered framework offers a powerful, scalable, and efficient solution to the contemporary problems in cybersecurity, and it has

high prospects of the practical application in real-time intrusion detectors.

## REFERENCES

- [1] Tavallae, M., Bagheri, E., Lu, W., & Ghorbani, A. A. (2009). "A detailed analysis of the KDD CUP 99 dataset." *Proceedings of the IEEE Symposium on Computational Intelligence for Security and Defense Applications*.
- [2] Sharafaldin, I., Lashkari, A. H., & Ghorbani, A. A. (2018). "Toward generating a new intrusion detection dataset and intrusion traffic characterization." *Proceedings of the 4th International Conference on Information Systems Security and Privacy (ICISSP)*.
- [3] Breiman, L. (2001). "Random Forests." *Machine Learning*, 45(1), 5–32.
- [4] Zhang, J., Zulkernine, M., & Haque, A. (2019). "Random-forest-based network intrusion detection systems." *IEEE Transactions on Systems, Man, and Cybernetics*.
- [5] Kim, G., Lee, S., & Kim, S. (2014). "A novel hybrid intrusion detection method integrating anomaly detection with misuse detection." *Computers & Security*, 41, 14–28.
- [6] Dua, S., & Du, X. (2020). *Data Mining and Machine Learning in Cybersecurity*. CRC Press.
- [7] Ferrag, M. A., Maglaras, L., Moschogiannis, S., & Janicke, H. (2020). "Deep learning for cyber security intrusion detection: Approaches, datasets, and comparative study." *Journal of Information Security and Applications*, 50, 102419.
- [8] Ahmad, Z., Shahid Khan, A., Shiang, C. W., Abdullah, J., & Ahmad, F. (2021). "Network intrusion detection system: A systematic study of machine learning and deep learning approaches." *IEEE Access*, 9, 124–150.
- [9] Sarker, I. H. (2021). "Machine learning for intelligent data analysis and automation in cybersecurity." *Annals of Data Science*, 8(4), 803–829.
- [10] Chandola, V., Banerjee, A., & Kumar, V. (2020). "Anomaly detection: A survey." *ACM Computing Surveys (CSUR)*, 41(3), 1–58.
- [11] Vinayakumar, R., Soman, K. P., & Poornachandran, P. (2019). "Applying deep

- learning approaches for network traffic prediction.” *IEEE Access*, 7, 143–156.
- [12] Aljawarneh, S., Aldwairi, M., & Yassein, M. B. (2018). “Anomaly-based intrusion detection system through feature selection analysis and building hybrid efficient model.” *Journal of Computational Science*, 25, 152–160.
- [13] Moustafa, N., & Slay, J. (2016). “The UNSW-NB15 dataset for network intrusion detection systems.” *Military Communications and Information Systems Conference (MilCIS)*.
- [14] Kasongo, S. M., & Sun, Y. (2020). “A deep learning method with wrapper-based feature extraction for wireless intrusion detection system.” *Computers & Security*, 92, 101752.
- [15] Zhou, Y., Cheng, G., Jiang, S., & Dai, M. (2020). “Building an efficient intrusion detection system based on feature selection and ensemble classifier.” *Computer Networks*, 174, 107247.
- [16] Chawla, N. V., et al. (2002). “SMOTE: Synthetic Minority Over-sampling Technique.” *Journal of Artificial Intelligence Research*.
- [17] Bommert, A., et al. (2020). “Benchmark for filter methods for feature selection in high-dimensional classification data.” *Computational Statistics & Data Analysis*.
- [18] Li, Y., Xia, J., Zhang, S., Yan, J., & Ai, X. (2021). “An efficient intrusion detection system based on support vector machines and gradually feature removal method.” *Expert Systems with Applications*.
- [19] Verma, A., & Ranga, V. (2022). “Machine learning based intrusion detection systems for IoT applications.” *Wireless Networks*.
- [20] Niyaz, Q., et al. (2016). “A deep learning approach for network intrusion detection system.” *EAI International Conference*.