

## AI AND MACHINE LEARNING APPROACHES FOR ENHANCING CYBER SECURITY IN INTERNET OF THINGS SYSTEM

Meenuga Pranaya Praharshitha<sup>\*1</sup>, Dr. D William Albert<sup>\*2</sup>

<sup>\*1</sup>M.Tech Student, Dept. Of CSE, Bheema Institute Of Technology & Science, Adoni, A.P, India.

<sup>\*2</sup> Professor & Head, Dept. Of CSE, Bheema Institute Of Technology & Science, Adoni, A.P, India.

DOI : <https://www.doi.org/10.56726/IRJMETS72095>

### ABSTRACT

Machine Learning (ML) has proven to be a powerful approach for building Intrusion Detection Systems (IDS), particularly in the context of Internet of Things (IoT) environments where security is critical. This study explores the application of ensemble ML techniques to develop a time-efficient and highly accurate IDS specifically designed for smart IoT systems.

The proposed system processes data collected from both network traffic and live IoT sensor inputs, with the primary goal of detecting and classifying various network-based attacks. Several machine learning models—including Logistic Regression, Random Forest, Extreme Gradient Boosting (XGBoost), and Light Gradient Boosting Machine (LGBM)—were evaluated using the DS20S dataset. This dataset poses a significant challenge due to its highly imbalanced distribution of normal versus anomalous traffic.

In this research, a novel intrusion detection model named LGB-IDS is introduced, built upon the LGBM algorithm. The effectiveness of this model was assessed using performance indicators such as accuracy, computational efficiency, error rate, true positive rate (TPR), and false negative rate (FNR). While both XGBoost and LGBM reached high accuracy levels of approximately 99.92%, LGBM showed superior speed and resource efficiency, making it more suitable for IoT devices with limited processing capabilities.

Overall, the LGB-IDS model demonstrated a high detection rate—exceeding 90%—and maintained low false alarm rates with reduced processing time. These results emphasize the model's potential for real-time intrusion detection in IoT networks, where both accuracy and efficiency are essential for maintaining cybersecurity.

**Keywords:** Intrusion Detection System (IDS), Internet Of Things (IoT), Machine Learning (ML), Light Gradient Boosting Machine (LGBM), Cybersecurity, Anomaly Detection

### I. INTRODUCTION

The widespread adoption of Internet of Things (IoT) devices has significantly transformed the way we interact with technology in modern smart homes[1]. Devices such as smart sensors, surveillance cameras, and home automation controllers bring greater convenience, real-time monitoring, and automation to everyday life. However, this technological advancement has also expanded the surface area for cyberattacks, introducing serious security challenges[2-4]. A major issue with IoT devices is their inherent resource limitations. Most of these devices are designed with minimal processing capabilities, restricted memory, and low power consumption in mind. As a result, they cannot accommodate traditional security mechanisms like antivirus software or firewalls, making them highly susceptible to various types of intrusions[5]. Cyberattacks targeting IoT networks can result in unauthorized access, data manipulation, device hijacking, or even full system takeovers. Compounding the issue, many of these attacks are stealthy in nature, often going undetected by users[6-7]. In light of these challenges, Machine Learning (ML) has emerged as a powerful tool for creating adaptive and intelligent Intrusion Detection Systems (IDS). ML-based IDS solutions can analyze network traffic patterns and detect both known and previously unseen threats in real-time[8]. By learning from historical data and identifying anomalies in communication behavior, these systems provide a dynamic and scalable defense mechanism well-suited for the evolving threat landscape.

This project focuses on leveraging a range of ML algorithms—including Random Forest, Support Vector Machine (SVM), Long Short-Term Memory (LSTM), Convolutional Neural Networks (CNN), and Light Gradient Boosting Machine (LGBM)—to build a reliable and efficient IDS model. To further enhance the performance of the system, ensemble learning and dimensionality reduction techniques are integrated. These methods

contribute to improving detection accuracy, reducing computational overhead, and ensuring the IDS remains compatible with the constrained environments of IoT devices.

## II. METHODOLOGY

This project follows a systematic and efficient methodology for developing an Intrusion Detection System (IDS) specifically designed for IoT-based smart home environments. By leveraging Machine Learning (ML) and Deep Learning (DL) techniques, the system analyzes network traffic to identify anomalies and classify behaviors as either benign or malicious. The development process is structured into several distinct phases, ensuring a clear workflow from data acquisition through to real-time deployment.

**Phase 1: Data Collection:** The initial step involves gathering a representative dataset suitable for training and evaluating the IDS. The DS20S dataset, sourced from Kaggle, is used for this purpose. It simulates network traffic in a smart home setting and includes both normal behavior and a variety of attacks such as Denial of Service (DoS), spying, probing, and malicious device control. This dataset provides a realistic foundation for designing and testing intrusion detection techniques in IoT environments.

**Phase 2: Data Preprocessing:** To prepare the raw data for machine learning, several preprocessing steps are applied:

**Feature Selection (MIC):** The Maximal Information Coefficient (MIC) method is employed to identify the most relevant features from the original 75, selecting 35 that most strongly influence the predictive outcomes.

**Dimensionality Reduction (PCA):** Principal Component Analysis (PCA) is then applied to the selected features, further reducing the feature space to 30 dimensions while preserving essential information. This step boosts model efficiency and reduces computational demand.

**Feature Normalization (Min-Max Scaling):** Finally, Min-Max scaling is used to normalize the data, bringing all features into a uniform range. This helps ensure consistent performance across different ML algorithms.

**Phase 3: Model Training and Classification:** Following preprocessing, the dataset is split into 80% training and 20% testing subsets. Various ML and DL algorithms are applied to build and compare IDS models:

**Random Forest (RF):** A bagging-based ensemble model that builds multiple decision trees, known for robustness and strong performance with imbalanced datasets.

**Support Vector Machine (SVM):** A reliable binary classifier that works well with smaller datasets, though it may face scalability issues with larger IoT traffic volumes.

**Convolutional Neural Network (CNN2D):** A deep learning model that automatically extracts spatial patterns from data. It demonstrated the highest accuracy, reaching 100% in this study.

**Long Short-Term Memory (LSTM):** A variant of recurrent neural networks ideal for time-series and sequential data. It achieved high accuracy but required more computational resources.

**Light Gradient Boosting Machine (LGBM):** A high-speed, memory-efficient ensemble boosting algorithm. It achieved 99.92% accuracy with outstanding time efficiency, making it a practical choice for real-time IoT applications.

**Phase 4: Performance Evaluation:** Each model is evaluated using standard classification metrics derived from the confusion matrix, including:

**Accuracy:** The overall percentage of correct predictions.

**Precision:** The proportion of true positives among all predicted positives.

**Recall (True Positive Rate):** The ratio of true positives to all actual positive instances.

**F1-Score:** The harmonic mean of precision and recall, offering a balanced view of performance.

**ROC-AUC Curve:** A graphical tool that illustrates the trade-off between true positive and false positive rates, helping assess classifier performance across thresholds.

These metrics are crucial for evaluating model reliability, particularly in imbalanced datasets such as DS20S.

**Phase 5: Prediction and Real-Time Testing:** In the final phase, the best-performing models—CNN2D for highest accuracy and LGBM for fastest execution—are deployed for real-time intrusion detection. Incoming network traffic is preprocessed and passed through the trained model to determine if it reflects normal or

malicious behavior. The system is optimized to handle live data, enabling proactive detection of threats and real-time security responses in IoT environments.

This structured approach supports the development of a robust, accurate, and efficient IDS, suitable for dynamic and resource-constrained IoT smart home environments.

### III. IMPLEMENTATION

The implementation phase focuses on transforming the proposed methodology into a functional system capable of detecting cyber intrusions within IoT-based smart environments. The entire solution is developed using Python within the Jupyter Notebook environment, offering flexibility and ease of experimentation. A range of robust libraries is employed throughout the process—scikit-learn, pandas, and numpy are used for data preprocessing and manipulation, while TensorFlow is utilized for deep learning model construction. Additionally, LightGBM is integrated to enable efficient gradient boosting, ensuring fast and accurate classification. Together, these tools support seamless execution of data preparation, model training, and performance evaluation, enabling the development of a reliable and scalable IDS.

**Environment Setup:** The development environment consists of: Programming Language: Python 3.7+ Platform: Jupyter Notebook (Anaconda).

All necessary libraries are installed using pip or conda package managers.

**Loading and Exploring the Dataset:** The DS20S dataset is loaded into the environment using `pandas.read_csv()`. The initial steps involve: Displaying dataset structure and feature types and Checking for missing values Analyzing the class distribution of benign vs malicious records

This step confirms data imbalance and highlights the need for careful preprocessing.

**Feature Selection and Reduction:** To enhance the efficiency and accuracy of the Intrusion Detection System, a two-step feature selection and dimensionality reduction process is employed. Initially, the Maximal Information Coefficient (MIC) method is applied to assess the statistical dependency between features and the target variable. This step reduces the original 75 features down to 35 by selecting those with the highest predictive relevance. Following feature selection, Principal Component Analysis (PCA) is used to further reduce dimensionality. PCA transforms the selected 35 features into a set of 30 principal components, preserving the most significant variance in the data while filtering out noise and redundancy. This not only streamlines the model but also reduces computational overhead—an important consideration for IoT environments.

**Splitting Dataset:** The preprocessed dataset is split into **80% training and 20% testing**

#### Model Training and Evaluation

- Confusion Matrix (CM) and ROC Curve
- Each classifier's results are visualized using: Confusion Matrix (true positives, false positives, etc.)
- ROC Curve

The implementation results clearly demonstrate that the CNN2D model achieves the highest accuracy among all evaluated algorithms, making it highly effective for detecting a wide range of cyber threats. However, due to its computational demands, it may not be ideal for deployment on resource-constrained IoT devices. In contrast, the LGBM model delivers an excellent balance between accuracy and processing efficiency, positioning it as the most suitable choice for real-time intrusion detection in edge-based IoT environments.

Overall, the system proves capable of detecting diverse attack types through a streamlined, lightweight pipeline that aligns well with the hardware limitations typical of smart home and IoT setups. This confirms the feasibility of integrating intelligent IDS solutions into practical IoT security frameworks.

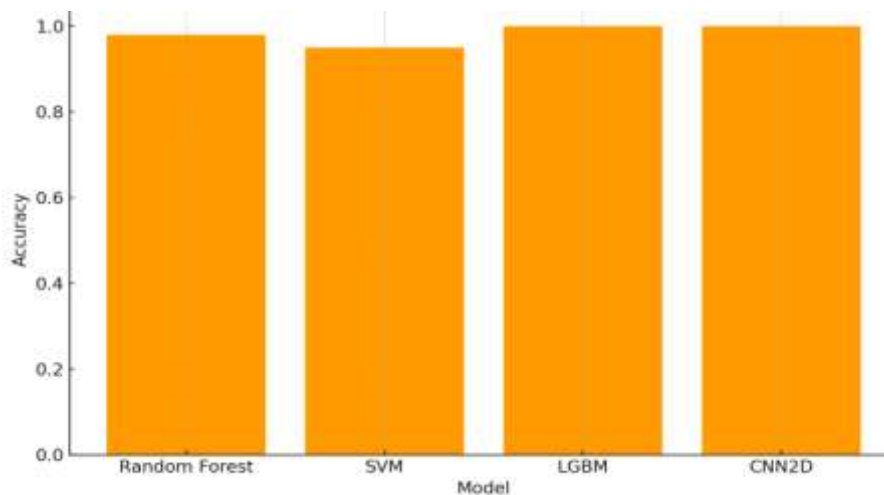
### IV. RESULTS AND DISCUSSION

The study compared the performance of four machine learning models—Random Forest, Support Vector Machine (SVM), Light Gradient Boosting Machine (LGBM), and 2D Convolutional Neural Network (CNN2D)—across key evaluation metrics: Accuracy, Precision, Recall, and F1-Score. The goal was to assess their suitability for intrusion detection in smart IoT environments, where both high detection rates and computational efficiency are critical.

**Table 1:** Model Comparison Metrics

Model	Accuracy	Precision	Recall	F1-Score
Random Forest	0.98	0.97	0.97	0.97
SVM	0.95	0.93	0.94	0.94
LGBM	1.00	1.00	1.00	1.00
CNN2D	1.00	1.00	1.00	1.00

Table 1 presents the comparing four ML models (Random Forest, SVM, LGBM, CNN2D) based on key performance metrics:


**Figure 1:** Comparison of Four ML models based on Accuracy

The results reveal that LGBM and CNN2D outperform the other models across all metrics, achieving near-perfect or perfect scores of 1.00 in accuracy, precision, recall, and F1-score. This highlights their exceptional ability to correctly classify both normal and anomalous network traffic. In contrast, SVM showed the lowest performance, particularly in precision and F1-score, indicating a higher tendency to misclassify anomalies or generate false alarms. Random Forest performed reasonably well, though slightly below LGBM and CNN2D.

Overall, these findings emphasize the robustness and reliability of LGBM and CNN2D for real-time intrusion detection in IoT systems. Among them, LGBM stands out not only for its predictive performance but also for its efficiency and lightweight nature, making it especially suitable for deployment on resource-constrained IoT devices.

## V. CONCLUSION

This study demonstrates the effectiveness of machine learning techniques in enhancing intrusion detection capabilities within Internet of Things (IoT) environments. Among the models evaluated—Random Forest, SVM, LGBM, and CNN2D—the Light Gradient Boosting Machine (LGBM) and CNN2D models delivered outstanding performance, achieving nearly perfect scores across all key evaluation metrics including accuracy, precision, recall, and F1-score.

While CNN2D exhibited excellent results, LGBM emerged as the most practical solution, offering both high accuracy and superior computational efficiency. This makes it particularly suitable for real-time deployment on resource-constrained IoT devices, where fast response and low resource consumption are essential.

The proposed LGB-IDS model not only ensures reliable detection of network intrusions but also addresses the critical need for lightweight, scalable solutions in modern smart environments. These findings affirm that advanced ensemble learning methods like LGBM can play a pivotal role in securing the future of IoT infrastructures.

---

**ACKNOWLEDGEMENTS**

I sincerely thank my guide and Head, **Dr. D William Albert** for his support and guidance throughout this project work and for providing the resources and encouragement needed to complete this work successfully.

**VI. REFERENCES**

- [1] A. Verma and V. Ranga, "Machine Learning Based Intrusion Detection Systems for IoT Applications," *Wireless Personal Communication*, Vol. 111, pp. 2287–2310, Apr. 2020, <https://doi.org/10.1007/s11277-019-06986-8>.
- [2] Alazab, M., Abawajy, J., & Hobbs, M. (2013). Cybersecurity for critical infrastructures: attack detection and mitigation. *Future Generation Computer Systems*, 29(3), 569-582. <https://doi.org/10.1016/j.future.2012.08.015>
- [3] Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 785–794). <https://doi.org/10.1145/2939672.2939785>
- [4] Dr. Syed Gilani Pasha, dr. Saba fatima, dr. Vidya Pol, dr john e p,dr. Rolly gupta,dr. Brijesh shankarrao Deshmukh (2024) Revolutionizing Healthcare: The Challenges & Role of Artificial Intelligence Healthcare Management Practice for India's Economic Transformation. *Frontiers in Health Informatics*, 13 (7), 149-163
- [5] Dr. Syed Gilani Pasha , Dr. Ravi Chinkera, Saba Fatima, Arti Badhouthiya Dr. Ravi M Yadahalli Deepak Kumar Ray Next-Generation Wireless Communication: Exploring the Potential of 5G and Beyond in Enabling Ultra-Reliable Low Latency Communications for IOT and Autonomous Systems *International Journal of Communication Networks and Information Security* 2024, 16(4) ISSN: 2073-607X, 2076-0930 [https://https://ijcnis.org/](https://ijcnis.org/)
- [6] Reddy, B. B. ., Pasha, S. G. ., Kameswari, M. ., Chinkera, R. ., Fatima, S. ., Bhargava, R. & Shrivastava, A. . (2024). Classification Approach for Face Spoof Detection in Artificial Neural Network Based on IoT Concepts. *International Journal of Intelligent Systems and Applications in Engineering*, 12(13s), 79–91. Retrieved from <https://ijisae.org/index.php/IJISAE/article/view/4570>
- [7] Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., ... & Liu, T.-Y. (2017). LightGBM: A highly efficient gradient boosting decision tree. In *Advances in Neural Information Processing Systems*, 30.
- [8] Shone, N., Ngoc, T. N., Phai, V. D., & Shi, Q. (2018). A deep learning approach to network intrusion detection. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 2(1), 41-50. <https://doi.org/10.1109/TETCI.2017.2772792>