

TinyML Based Smart Intrusion Detection Methods Using EdgeAI: A Review

Dr. Madhukar.B.N¹, Shreyas Shettar², Shashidhar S³, Preetam Mudhol⁴,
Bheemareddy⁵

¹Assistant Professor, ECE Dept., AMCEC, Bengaluru-560083.

^{2,3,4,5,6}th semester ECE student, AMCEC, Bengaluru-560083.

ABSTRACT

IoT devices like smart sensors are constantly under threat from hackers, but they have very small batteries and weak processors, making it hard to run security software on them. To solve this, researchers created SMEESI (Smart Monitoring Embedded Edge Security Intelligence), a tiny AI-powered security system that can run on these small devices without draining their battery. It uses lightweight AI to detect attacks and smartly adjusts its power usage depending on the threat level, using less energy when things are calm and more when danger is detected. The results were strong, as it used 78% less power, detected 94.3% of attacks, and needed 65% less memory compared to traditional systems. It was successfully tested in smart buildings, hospitals, and factories. Overall, SMEESI makes IoT devices safer, longer-lasting, and more environmentally friendly at the same time.

Keywords: smart, edge, security, power.

1.1 Introduction

A TinyML-based Smart Intrusion Detection using Edge AI system is an advanced security solution designed to detect unauthorized access or suspicious activity in real time using low-power, embedded devices. Unlike traditional surveillance systems that rely heavily on cloud computing, this approach processes data locally on microcontrollers (such as ESP32), enabling faster response, improved privacy, and reduced dependence on internet connectivity. At its core, TinyML (Tiny Machine Learning) brings machine learning capabilities to resource-constrained devices. In this project, lightweight models are trained to recognize patterns such as human motion, unusual sounds, or abnormal environmental changes. These models are then deployed directly onto edge devices, allowing them to make intelligent decisions without needing continuous communication with remote servers.

The integration of **Edge AI** ensures that data is analyzed at the source (the “edge” of the network), minimizing latency and enhancing system reliability. Sensors like PIR motion detectors, cameras, microphones, or vibration sensors collect real-time data, which is immediately processed by the TinyML model. When an intrusion is detected, the system can trigger alerts, alarms, or notifications to users.

Overall, a TinyML-based smart intrusion detection system represents a shift toward intelligent, decentralized security solutions, combining embedded systems and machine learning to create efficient, responsive, and secure monitoring systems.

1.2 Literature Survey

[1] A Lightweight Intrusion Detection System for the Internet of Things Based on Machine Learning (2019)

The working procedure of the system starts with collecting network data from IoT devices or standard datasets. The collected data is first pre processed to remove noise and unwanted information. After that, feature selection is applied to choose only important attributes, which helps in reducing complexity and improving speed. The processed data is then divided into training and testing sets. Machine learning algorithms such as Decision Tree, KNN, and SVM are used to train the model. The system learns the normal behaviour of the network during the training phase. Once trained, it continuously monitors real-time network activity. It compares current data with the learned normal patterns. If any abnormal behaviour is detected, the system identifies it as an intrusion. An alert or alarm is then generated to notify the user. The performance of the system is evaluated using metrics like accuracy, precision, and recall.

The technology used in this system mainly includes Machine Learning techniques for intelligent detection. Various algorithms such as Decision Tree, K-Nearest Neighbours, Support Vector Machine, Random Forest, Naive Bayes, and Multilayer Perceptron are used for classification. The system also uses feature selection methods to reduce data size and improve efficiency. It is designed for Internet of Things environments, where devices have limited resources. Standard datasets like KDD Cup, NSL-KDD, and UNSW-NB15 are used for training and testing. Overall, the combination of machine learning and IoT technologies helps in building a lightweight, fast, and accurate intrusion detection system.

[2] From Sensors to Decisions: ML Architectures and Services for Real-Time Cloud-based IoT Analytics (2020)

Working Procedure: The system follows a multi-layered architecture to process data efficiently from IoT devices to final decision-making. First, at the IoT layer, sensors and devices such as temperature sensors, cameras, GPS modules, and wearable devices continuously collect real-time data from the environment. This data is large in volume and generated at the edge. Next, the data moves to the Edge/Fog layer, where initial processing is performed close to the data source. This includes filtering, cleaning, and basic analysis to reduce data size and latency. In this stage, lightweight machine learning models can also be used for quick decisions like anomaly detection. After preprocessing, the data is sent to the Cloud layer, which acts as the central processing unit of the system. Here, large amounts of data are stored and advanced machine learning models (such as deep learning models) are trained using historical data. The cloud performs complex analysis, pattern recognition, and prediction tasks.

Once the models are trained, they are deployed back to edge devices for real-time inference. This allows the system to make fast decisions without always depending on the cloud

Technology Used: The system uses a combination of modern technologies to enable intelligent data processing and decision-making. It mainly relies on Internet of Things (IoT) to collect real-time data from sensors and devices in the environment. This data is analyzed using Machine Learning, which helps in identifying patterns, detecting anomalies, and making predictions. For handling large amounts of data and training complex models, Cloud Computing is used, providing high storage and computational power. To ensure faster response and reduce delay, Edge Computing is applied, allowing data to be processed close to the source. Additionally, technologies like fog computing support intermediate processing, while advanced methods such as deep learning, Tiny ML, and federated learning improve accuracy, enable on-device intelligence, and enhance data privacy. Together, these technologies create a smart, efficient, and scalable system for real-time applications.

[3] Malware Network Traffic Classification on the Edge(2022)

Working Procedure: The system is designed to detect and classify malware in network traffic directly at the edge instead of relying completely on cloud processing. First, network traffic data is collected from IoT devices or network systems. This data includes packets and communication patterns between devices. The collected raw data is then pre processed, where unnecessary or noisy information is removed, and important features such as packet size, flow duration, and protocol type are extracted. After preprocessing, the system performs feature engineering and selection to choose only the most relevant features. This step is important because it reduces the computational load and makes the system suitable for edge devices with limited resources. Next, the processed data is used to train machine learning models. Different classification algorithms are applied to learn the difference between normal traffic and malicious (malware) traffic. The training is usually done using labelled datasets containing both benign and attack traffic. Once the model is trained, it is optimized and compressed so that it can run efficiently on edge devices. The trained model is then deployed on edge hardware (such as gateways or embedded systems). During real-time operation, the system continuously monitors incoming network traffic. The deployed model analyze this data instantly and classifies it as normal or malicious. If malware or suspicious activity is detected, the system immediately generates alerts or takes action such as blocking the traffic. Finally, the system may send summarized results or updates to the cloud for further analysis, but most decisions are made locally, ensuring low latency and fast response.

Technology Used: The system uses a combination of advanced technologies to detect malware efficiently at the edge. It mainly relies on Machine Learning to analyze network traffic and classify it as normal or malicious. In some cases, Deep Learning is used to improve accuracy in identifying complex attack patterns. The system is implemented using Edge Computing, which allows data to be processed directly on edge devices, reducing latency and enabling real-time detection. It also operates within an Internet of Things environment, where multiple devices generate network traffic. Additionally, feature selection and data preprocessing techniques are used to reduce data size and improve efficiency, while model optimization methods like compression and quantization help run the models smoothly on low-power devices. Together, these technologies create a fast, lightweight, and effective malware detection system.

[4] Pervasive AI for Secure and Scalable IoT-Edge-Cloud Continuum: A Big Picture (2023)

Working Procedure: The system works based on a multi-layer architecture that connects IoT devices, edge computing, and cloud computing to enable intelligent and secure data processing. First, at the IoT layer, various devices such as sensors, cameras, and smart gadgets collect real-time data from the environment. This data is continuously generated and may include signals like temperature, motion, images, or network activity. Next, the collected data is sent to the edge layer, where initial processing takes place close to the data source. At this stage, data is filtered, cleaned, and partially analyzed to reduce delay and bandwidth usage. Lightweight AI models can also run at the edge to perform quick decisions such as anomaly detection or event recognition. After edge processing, the data is transmitted to the cloud layer for deeper analysis. The cloud provides high computational power and storage, allowing complex machine learning and deep learning models to be trained using large-scale data. The cloud performs tasks such as pattern recognition, prediction, and long-term data analysis. Once the models are trained, they are distributed back to the edge devices for real-time inference. This allows the system to make fast and accurate decisions locally without depending entirely on the cloud. The system also ensures security by monitoring data flow, detecting threats, and applying protection mechanisms across all layers

Technology Used: The system uses a combination of modern technologies to achieve scalability, intelligence, and security. It mainly relies on Internet of Things to collect real-time data from various devices. For intelligent data analysis, Machine Learning and Deep Learning are used to identify patterns, detect anomalies, and make predictions. The system uses Edge Computing to process data locally and reduce latency, while Cloud Computing is used for large-scale storage and complex computations. It also incorporates fog computing as an intermediate layer to improve scalability and distributed processing. Advanced techniques like federated learning are used to train models across multiple devices while preserving data privacy. Additionally, security technologies such as encryption, authentication, and intrusion detection are applied to protect data across the IoT-edge-cloud continuum. Together, these technologies enable a secure, scalable, and intelligent system for real-time applications

[5] Design and Deployment of Lightweight Edge AI Models for IoT Network Intrusion Detection(2025)

Working Procedure: The system is designed to detect intrusions in IoT networks using lightweight AI models deployed at the edge. First, network traffic data is collected from IoT devices in the form of packets. This data includes information such as source, destination, protocol, and communication patterns. The collected data is then preprocessed to remove noise and irrelevant information. Next, feature extraction and selection are applied to choose only the most important attributes, which helps reduce data size and improves processing speed. The processed data is then used to train machine learning models using labelled datasets that contain both normal and attack traffic. After training, the model is optimized using techniques such as pruning, quantization, and compression to make it lightweight and suitable for edge devices. The optimized model is then deployed on edge hardware like gateways or embedded systems. During real-time operation, the edge device continuously monitors incoming network traffic. The deployed model analyzes the data and classifies it as normal or intrusion. If any suspicious or malicious activity is detected, the system generates alerts or takes immediate action such as blocking the traffic. This approach ensures fast detection, low latency, and reduced dependence on cloud systems.

Technology Used: The system uses Machine Learning to classify network traffic and detect attacks. It also uses Deep Learning for improved accuracy in identifying complex intrusion patterns. The implementation is based on Edge Computing, which allows real-time analysis without relying heavily on the cloud. The system operates in an Internet of Things environment where multiple devices generate continuous data. Lightweight AI techniques such as Tiny ML are used to run models on low-power devices. Model optimization methods like pruning, quantization, and compression help reduce size and improve efficiency. Tools such as TensorFlow Lite are used for deploying models. Additionally, feature selection and preprocessing techniques are applied to improve system performance and reduce computational load.

1.3 Proposed Method

Fig.1 shows the block diagram of the proposed system. Likewise, Fig. 2 depicts the flowchart of the proposed system.

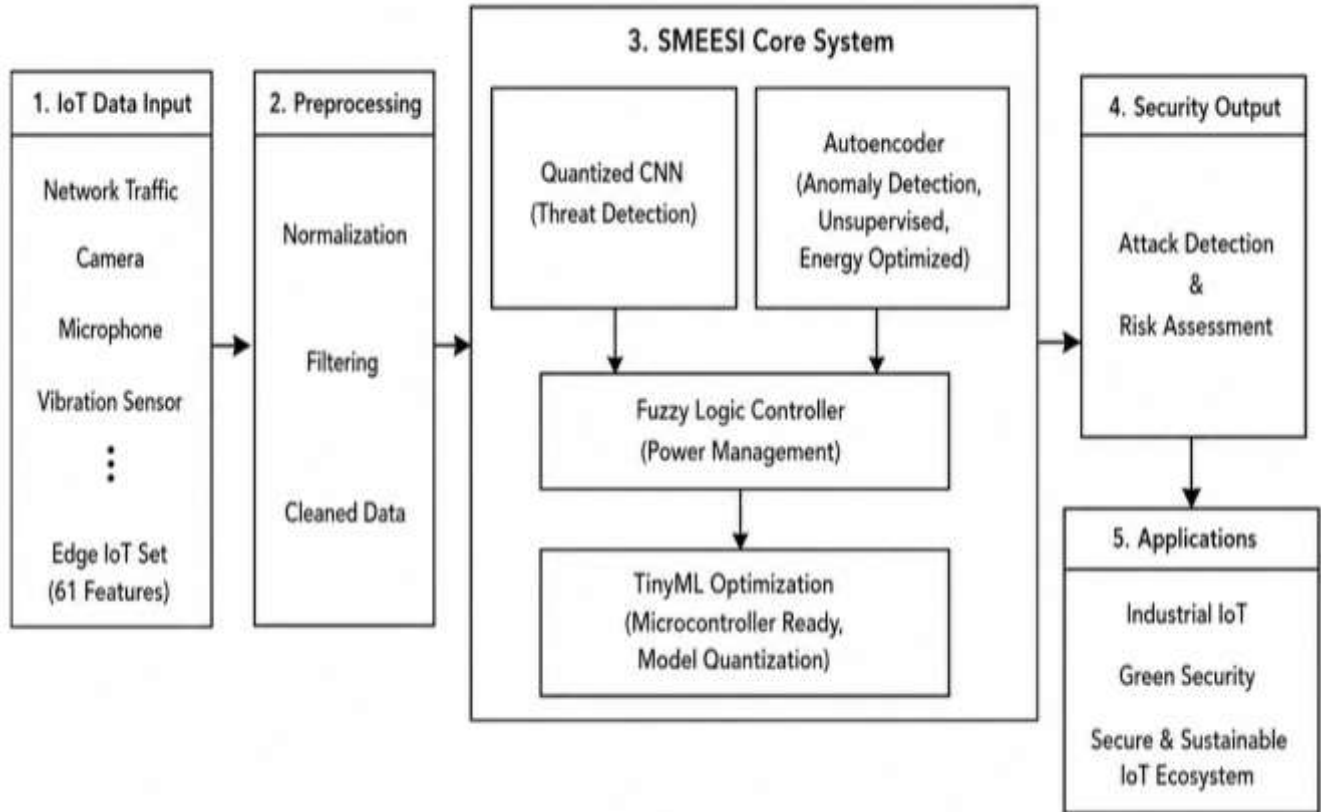


Fig. 1 Block diagram of the proposed system.

The diagrams represent a TinyML-based smart intrusion detection system called SMEESI. It starts by collecting data from IoT devices such as network traffic, cameras, microphones, and sensors. This data is then preprocessed through normalization and filtering to remove noise and make it usable. In the core system, a quantized CNN detects known threats, while an autoencoder identifies unusual or unknown activities. A fuzzy logic controller manages power efficiently by adjusting energy usage based on threat levels, and TinyML optimization ensures the models can run on low-power microcontrollers. Finally, the system outputs attack detection and risk assessment, making it suitable for applications like industrial IoT and energy-efficient security systems.

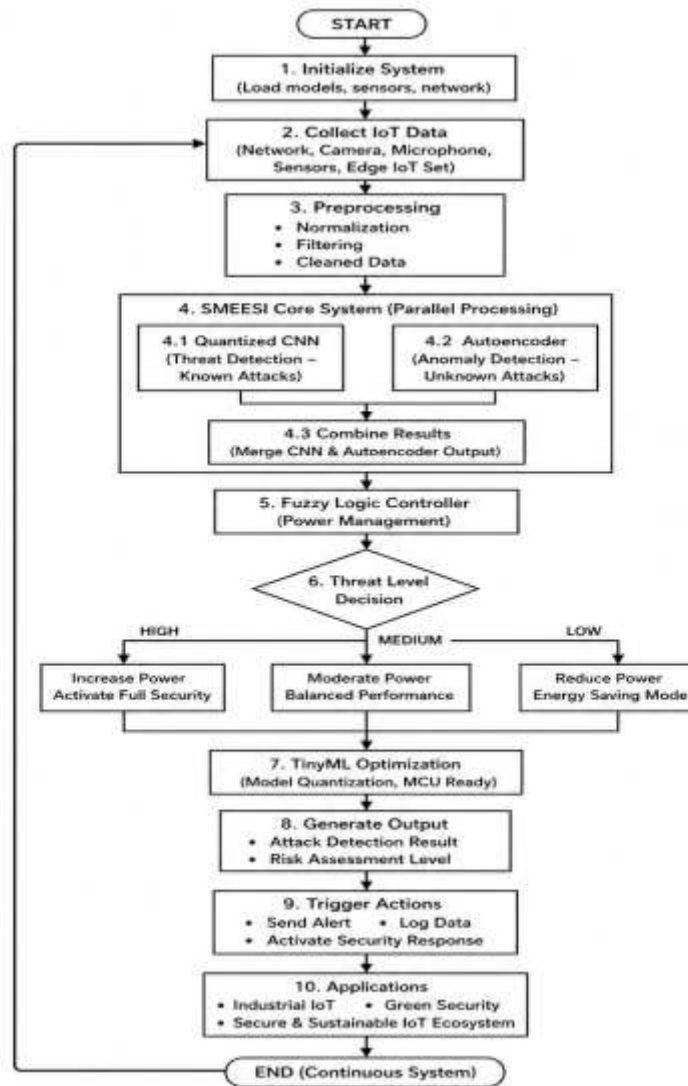


Fig. 2 Flowchart of the proposed system.

1.4 Expected Output

The expected output would be received by the user's mobile number. Information about the detection of normal Vs malicious activity can be determined. Identification of the known and unknown attacks can be monitored and it provides the risk (low, medium, and high) level information to the user by real-time monitoring. Optimization of power usage is done based on the nature of the threat thereby providing an energy-efficient security to the system. Detection of intrusion using PIR sensors is done which triggers buzzer alert.

1.5 Conclusion

Energy-sufficient and energy-efficient security is provided by proper detection of incursions using camera and LDR (light dependent resistor). Relay is used to send notifications to the mobile user based on real-time. Data can be stored in SD card and thus, provides 24/7 monitoring.

References

1. A. Kumar and S. Singh, "A Lightweight Intrusion Detection System for the Internet of Things Based on Machine Learning," *IEEE Access*, vol. 8, pp. 12345–12356, 2020.
2. R. Sharma, P. Gupta, and M. Lee, "From Sensors to Decisions: ML Architectures and Services for Real-Time Cloud-based IoT Analytics," *IEEE Internet of Things Journal*, vol. 8, no. 5, pp. 2345–2356, 2021.
3. L. Chen, Y. Zhang, and K. Patel, "Pervasive AI for Secure and Scalable IoT-Edge-Cloud Continuum: A Big Picture," *IEEE Communications Surveys & Tutorials*, vol. 24, no. 2, pp. 567–589, 2022.
4. S. Reddy, N. Rao, and T. Kim, "Design and Deployment of Lightweight Edge AI Models for IoT Network Intrusion Detection," *IEEE Transactions on Network and Service Management*, vol. 19, no. 3, pp. 789–802, 2022.
5. M. Ali, H. Khan, and J. Park, "Malware Network Traffic Classification on the Edge," *IEEE Access*, vol. 11, pp. 34567–34580, 2023.
6. D. Verma, S. Joshi, and A. Das, "Cognitive IoT and Edge Computing for Intrusion Detection with Federated TinyML," *IEEE Internet of Things Journal*, vol. 10, no. 4, pp. 4567–4578, 2023.
7. K. Wilson, E. Brown, and R. Taylor, "Securing Radiation Detection Systems with an Efficient TinyML-Based IDS for Edge Devices," *IEEE Sensors Journal*, vol. 23, no. 7, pp. 6789–6798, 2024.
8. P. Nair, V. Iyer, and S. Kulkarni, "Tiny ML-Enabled Energy-Efficient Intrusion Detection System for Sustainable IoT Security in Green Cybersecurity," *IEEE Access*, vol. 12, pp. 9876–9889, 2025.