

# In-Depth Analysis of Machine Learning for Securing Internet of Things devices using CIC IoT & Net Flow Dataset

MSc Research Project  
Cybersecurity

Shreyas Srinivasa  
Student ID: X23102641

School of Computing  
National College of Ireland

Supervisor: Dr. Arghir Moldovan

**National College of Ireland**  
**MSc Project Submission Sheet**  
**School of Computing**



**Student Name:** Shreyas Srinivasa  
**Student ID:** 23102641  
**Programme:** Cybersecurity **Year:** 2024-2025  
**Module:** MSc Research Project  
**Supervisor:** Dr. Arghir Moldovan  
**Submission Due Date:** 12-12-2024  
**Project Title:** In-Depth Analysis of Machine Learning for Securing Internet of Things devices using CIC IoT and Net Flow Dataset  
8282  
**Word Count:** ..... **Page Count:**.....18.....

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

**Signature:** Shreyas Srinivasa  
12-12-2024  
**Date:** .....

**PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST**

|   |                          |
|---|--------------------------|
| Attach a completed copy of this sheet to each project (including multiple copies)   | <input type="checkbox"/> |
| <b>Attach a Moodle submission receipt of the online project submission,</b> to each project (including multiple copies).  | <input type="checkbox"/> |
| <b>You must ensure that you retain a HARD COPY of the project,</b> both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer. | <input type="checkbox"/> |

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

| <b>Office Use Only</b>           |  |
|----------------------------------|--|
| Signature:                       |  |
| Date:                            |  |
| Penalty Applied (if applicable): |  |

# In-Depth Analysis of Machine Learning for Securing Internet of Things devices using CIC IoT and Net Flow Dataset

Shreyas Srinivasa  
X23102641

## Abstract

The rise of Internet of Things (IoT) has made the digital landscape transformed into providing more services but also introduced significant cybersecurity challenges by expanding potential vulnerabilities. Security systems such as Data Loss Prevention (DLP), Intrusion Detection Systems (IDS) and firewalls are struggling to keep up with these modern threats. They often produce a high number of false positives and lack the capability to identify more advanced, evolving attacks. To address these shortcomings, our research presents a machine learning models where we compare the ML models and analyse in detail as to which model produces the best accuracy. This approach has been documented each model and made sure the results are best for the networks for processing sequential data. The hybrid design improves detection accuracy and reduces the rate of false positives.

Furthermore, we developed these machine learning models are two datasets in order to figure out the best results. The novelty of this research lies in the hyper parameter tuning of machine learning models to achieve best results. The main contributions of this research are the extensive machine learning and deep learning algorithms for detecting and classifying malicious traffic. Our research conducted five novel machine learning models on UQ-NIDS dataset and added hyper parameter tuning in order to determine how well the model performs on those IoT attacks and for CICIoT dataset, our research is performed extensively on the classification of attack categories and implemented three novel algorithms to find the best model for each attack categories. And our research was able to find that Random Forest and Logistic Regression performs well on both the datasets and also identifying different categories in IoT attacks.

## 1 Introduction

The rapid raise in use of technology, the main focus of this is the security aspect as it plays a major role in securing the infrastructure or privacy of the user. There are many organizations which are using IoT are bent over backwards to keep up with ever evolving cybersecurity landscape. The protection of sensitive data and resources is not an option but rather a necessity. And without securing these important assets organizations are exposed to exploiting these vulnerabilities. The challenges faced by the cybersecurity professionals in overcome these challenges by developing counterattack strategies such DDoS (Distributed Denial of Service), Phishing attacks, Brute force attacks, Ransomware attacks. The risks of compromise have become higher and mitigating cyber security attacks is necessity and not a choice. All organizations should prioritize the implementation of the latest day modern tools and monitoring tools to always stay alert of the rapidly increasing cyber threats.

The addition of IoT devices to the current modern day world adds more complication to the issue and these companies should stay ahead of the the curve. The cybersecurity experts work on developing a proactive approach to overcome the cyber attacks and also to address the weakness in their services. Provided the fast growth of IoT devices threats are at all time highest making sophisticated and robust strategies as crucial aspect. Organizations should implement latest intrusion detection systems and monitoring services to fight against cyber threats and provide security to the resources. Additionally employee training are vital and business should implement these practices. The increasing cyber threats and risks associated with use of IoT devices used for commercial and household services emphasize the importance of effective threat detection to identify malicious intrusions.

The main challenge with network attacks for IoT devices is for the IRT (Incident Response Team) which has to be the strongest for providing any defence as they need to have a solid blueprint. With the ever changing and ever evolving network attacks and posing a great difficulty for the cyber security experts and consultants for easing any threats in the future. These attacks have to be studied and understood in order to detect and find new solutions to the cyber threats. And the new cyber attacks cannot be predicted where it might the effect the system. Countless measures and methods have been introducing in order to counter attack the existing challenges but the recent advancements in IoT devices require more advanced threat detection techniques to solidify the organization's property and privacy. There are many research done on understanding the impact of the best settings to mitigate the attacks and results captured for the intrusion detection in IoT environments. There were plenty of methods added but the network-based intrusion detection system and the machine learning based detection system were the latest technology that can be used to combat the cyber threats. The use of AI has become a top priority in any organization for intrusion detection as they are fast and operate without any human intervention.

The threats and cyber intrusion detection methods associated with the use of IoT devices in both household and commercial purpose, there is a massive need to develop effective cyber threat detection methods. And the most common as we all know is the intrusion detection system (IDS) but there are few security attacks such as DoS (Denial of Service) attacks, XSS (Cross Site Scripting) attacks which pose significant challenges, and it requires detection and prevention on how to identify the most the posed network attacks. Intrusion Detection Systems is effective for identifying the most cyber attack and these raise the configuration to be more sophisticated. Bakro et al (2023) emphasized the importance of introducing a secure IDS system which can be capable of detecting attacks. The IDS driven with AI technology should be able to any attack that pose a great threat to the system. And developing such integrated technology takes a lot in depth research and vibrant methodology to be introduced. The testing process will also imply if the AI driven IDS provide the best results and if it should be implemented.

As the organizations are increasingly upgrading their network, there is an urgent need for an effective threat detection mechanism which can make sure there is an incident response. All the artifacts suggest if integrating artificial intelligence into intrusion detection systems which

can enable precise and swift detection. The modern-day intrusion detection methods including deep learning use datasets to detect intrusions in an IoT device. Studies also emphasized the importance of AI for successful intrusion detection as it provides detailed insights for its solutions and along with its datasets. In this context both Net flow V1 and CiC IoT has significant attention because for its effectiveness in the domain.

The novelty of this research is tuning of machine learning models in order to achieve best results. The contributions of this research are:

- Deep learning algorithms for detecting and classifying malicious traffic.
- Conducted five novel machine learning models on UQ-NIDS dataset and added hyper parameter tuning.

## **1.1 Research Objectives**

In order to improve the IoT security and enhance it, we will these objectives.

- Evaluate the performance of Logistic Regression, Decision Tree, Support Vector Machine, Random Forrest, XG Boost and Naïve Bayes on both Netflow (UQ-NIDS) and CiC IoT dataset using metric such as accuracy, precision, recall and F1 score.
- Find out which of these listed machine learning models is the most suitable for anomaly detection in IoT data for the UQ-NIDS dataset and we have implemented hyperparameter tuning in order to enhance the received output.
- For the CIC IoT dataset the aim of contributing to the proposals of an efficient classification of these attacks into 34 class (33+1), 8 class (7+1) and 2 class (1+1) classification.

## **1.2 Research Question**

In this work the aim is to address the following question

- Which machine learning model should be applied in appropriate manner in IoT devices to detect the different attacks and how well these models function with hyper parameter tuning.

## **1.3 Contributions**

The contributions one can expect from our research is that which machine learning or deep learning model is suitable for IoT cyber threats. The in-depth analysis of the ML models on two unique datasets is done in order to figure out the best performing model and the attack categories in dataset are further classified in order to understand the attack type and how well the model performs against those attacks.

# **2 Related Work**

The study focuses more on the information from the literature which has been the important in the domain of network intrusion detection system (NIDS) and differs with the rest of the

methods in investigating cyber-attacks. The information has been contextualized for different datasets which are often used in detection process, and this gives the opportunity to explore the significance of these datasets for detecting novel attacks.

## **2.1 Techniques for Intrusion Detection and Cybersecurity Challenges**

This section is to provide discussion to understand and explore the gaps to identify the need for integrated intrusion detection as it is crucial due to the rising cyber threats in the network domain and applications. Guezzaz et al. (2022) has pointed that malicious and novel attacks has continued to pose great challenge for users which in turn has affected both wire and wireless networks. Similarly, Kayode Saheed et al. (2022) has discussed about the network threats in the IoT environments and noted that use of network application has directly increased the privacy threats as the emergence of novel attacks. Organizations all over the world are increasing their investment and also making efforts in order to enhance the detection of malicious attack, smart procedures are being examined by measuring accuracy with other parameters and by comparing them is to identify the most effective methods for recognizing attacks. Automated intrusion detection is vital due to the ever-evolving cybersecurity threats. These challenges impose great threat to the consumers that is the users as it affects the medium of communication.

Nizamudeen (2023) has highlighted that the intelligent classification of detection methods had garnered significant attention due to the development of advanced network systems and with the persistent threat of cyber-attacks. The assessment of this study is very important as it revealed Intelligent Intrusion Detection Framework is very crucial in identifying the novel network and application attacks. This framework makes use of various datasets to preprocess the information and select important features. Intrusion detection based on deep learning was known as 2D-ACNN and served as binary classifier for detecting both normal and abnormal network with accuracy rate more than 96%. Douiba et al. (2023) had discussed the significance of an IoT as a enhanced network infrastructure and also acknowledged the inevitable security threats that poses. The study also examined various machine learning and deep learning methods and achieved nearly 96.5% detection accuracy with less computation time. Network Intrusion Detection systems have become so important due to their enhanced defence mechanism against all these diverse attacks (Layeghy et al., 2024) and the study identified issues with datasets like KDD99 which have limitations which results in suboptimal outcomes. Thre study also recommends use of high-quality NIDS datasets.

Singh (2022) explored the fundamental concept of intrusion detection which highlights as the crucial network infrastructure. It also emphasized that malicious activities by users can create more vulnerabilities within the system. The study also focused on new intrusion detection methods examining key areas: analysis of IDS techniques, providing specifi insights and identifying the future research opportunities. The importance of detecting cybersecurity threats in IoT environments is because of the increasing vulnerabilities that are faced by users Olabanji et al. (2024). Their study also discussed the impact of artificial intelligence and various other algorithms in detecting the malicious attacks within network system and these systems provide hybrid security infrastructure with improved predictive capabilities and has also demonstrated enhanced accuracy in threat detection.

Cybersecurity is the most critical aspect in this modern-day technology and there are more research being conducted because of its increasing network and the system vulnerabilities.

The growing importance of artificial intelligence and machine learning based intrusion detection through the behaviour analysis as highlighted by Talukder et al.(2024) given the high detection accuracy with also providing dynamic cyber threats. The study also highlights the importance of datasets such as CIC-IDS 2017 and CIC-IDS-2018.

The machine learning models such as decision tree, random forests and ensemble classifiers achieved impressive results with 98% using CIC-IDS-2017 dataset. The other study by Zarpulao et al. (2017) emphasized IoT environment and addressed the network security issues, and the study was able to find out intrusion detection systems methods are ineffective in detecting novel cyber threats. And suggested the need for new IDS methods and explore further schemes to enhance IDS.

## **2.2 Innovation in Intrusion Detection System**

The analysis in recent growth trends in IT highlights both new opportunities and the worries (Archana et al 2021). Moreover Bharati & Tamane (2020) have also explored the effectiveness of using advanced NIDS methods which are driven by machine learning (ML) and deep learning in detecting the cyber-attacks. Long et al. (2024) has placed more importance on the need for more robust defence mechanism systems within the infrastructure in order to secure the network accessible resources. They also introduced transformer-based NIDS method that has achieved 95% accuracy rate and outperformed the other models which has effectively enhanced the security. AL-Ghuwairi et al. (2023) had noted the gradual increase in importance of NIDS for addressing security challenges. But the main issue was high false-positive alarm rate had persisted primarily due to limitation of existing datasets and to tackle this a new improved feature selection method was introduced that can address time series anomalies and the attacks. The results concluded by stating that feature selection prediction model had demonstrated enhanced performance. Attou, Guezzaz, et al. (2023) had explained that detecting anomalies in the network traffic requires a improved IDS system so has to ensure accurate predictions. The existing models fail to detect the novel attacks effectively and the new models have been developed.

Our literature review highlights several faults in intrusion detection systems (IDS). Many of the modern IDS models including the machine learning and deep learning suffer from high false positive rates and also struggle to detect novel cyber attacks in IoT environments. This research proposes modern IDS that is subjected to different machine learning models and have three different classifications to enhance detection accuracy and reduce false positives.

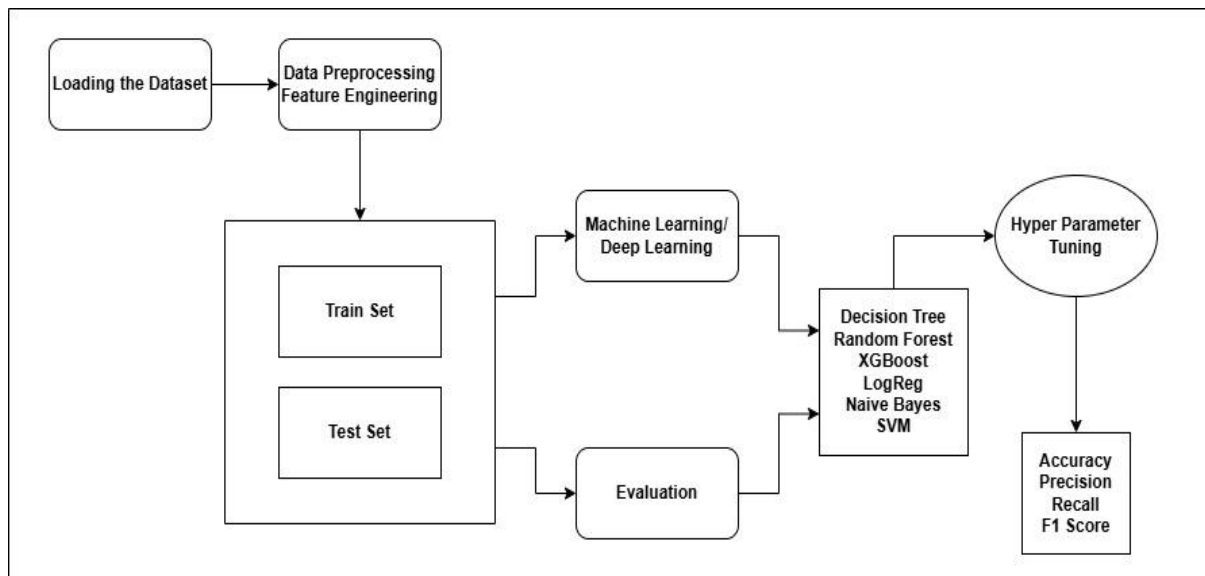
Nadia Chaabouni et al. (2019) provide vital information as to why to secure IoT networks and increasing number of cyber threats. This paper by Nadia focuses more Network Intrusion Detection systems which are tailored for IoT and emphasize more on using machine learning techniques. The document also reviews approaches to IoT NIDS and the limitations of traditional detection systems in IoT. It also presents an in-depth classification of IoT threats based on the architectural layers and design. The paper also evaluates existing datasets such as KDD99, NSL-KDD which have been benchmarks for training and testing intrusion detection systems. The paper also emphasizes the shift to hybrid detection techniques that utilize machine learning to enhance detection accuracy and reduce false positives which are very important for detecting threats and zero-day attacks. The paper underscores the underlying growth on machine learning including deep learning models to address IoT dynamic and more complex threats.

Marwa et al. (2023) explored the application of machine learning techniques for Intrusion Detection systems within IoT environments which have become increasingly vulnerable to security threats. The paper particularly emphasizes on enhancing IDS for handling unique challenges of IoT networks such as diverse communication and constrained computational resources. The study reviews existing machine learning techniques applied in IDS and commonly used algorithms such as Decision Tree, Support Vector Machines, Naïve Bayes and Random Forest, the authors conduct experimental comparisons using NSL-KDD dataset for binary and multiclass classification tasks and the dataset categorizes network activities to normal and various attack types. The authors affirm that machine learning offers significant potential for securing IDS to IoT networks.

Bambang Susilo et al. (2020) focuses on enhancing the IoT security through machine learning and deep learning techniques. IoT eco system connects billions of devices globally and faces significant security challenges because of its open nature. Traditional security systems like encryption and authentication fail especially against advanced threats like botnet and the study proposes using machine learning and deep learning to enhance IDS capabilities for IoT networks. The authors propose future framework work on hybrid models like combining machine learning and deep learning techniques to improve the detection capabilities and the paper provided a comprehensive analysis of IoT for machine learning and deep learning security which offered valuable insight insights.

### 3 Research Methodology

Our methodology includes several steps in order to develop an optimal intrusion detection system which is capable of detecting anomalies. The methodology diagram for both the datasets illustrates the research process below. To make sure that there is high level security is crucial for maintaining both safety and trustworthy communication between multiple organizations and this research also involves gathering raw data from various sources and preprocess it to identify a large number of unique values in the target and leads to multiclass classification. The methodology for the CICIoT 2023 is inspired from the exponential growth and integration of IoT technologies which introduce new challenges in terms of security. This dataset includes 105 real IoT devices from various categories. CICIoT dataset features 33 distinct attacks and includes 105 real IoT real devices from various categories.





### 3.1 Dataset Description

In this research, we are using two datasets provided by the University of Queensland (UQ-NIDS) and CIC IoT dataset. The UQ-NIDS consists of small files containing the network flows from various networks. The attack categories were modified by combining parent categories. Several attacks were grouped into brute force category including FTP, SSH and Brute Force-Web while SQL injection attacks were categorized under the injection attacks. The dataset comprises of 11,994,893 records with 9,208,048 (76.77%) classified as benign network flows and the remaining include 763,285 DDoS attacks, 482,946 reconnaissance attacks, 468,575 injection attacks, 348,962 DoS attacks, 291,955 brute-force attacks and 156,299 password attacks. Classes such as XSS, infiltration, exploits, scanning, fuzzers, backdoor, bot, generic, analysis, theft, shellcode, MITM, worms and ransomware are also present.

The CICIoT dataset contains over 548 GB of raw network traffic data which are stored in PCAP format using Wireshark. In order to make this dataset usable for machine learning 47 key features are extracted from the raw data. This dataset is available in two formats, the original PCAP files for advanced analysis and CSV files which contains extracted features for machine learning tasks. Malicious traffic has 33 different attack types which are categorized into 7 primary groups. There is a total of 46,686,579 records. DDoS category contains 33,984,560 records, DoS category comprises of 8,090,738 records which focuses on single source attacks which can disrupt the services availability. PSHACK flood comprises of 4,094,755, HTTP flood comprises of 28,790 records. SYN flood comprises of 2,028,834 records. The spoofing category has 486,504 records. And Mirai Botnet category consists of 2,634,124 records, Benign Traffic has 1,098,195 records. This comprehensive breakdown shows the diversity of CIC IoT dataset which makes it an invaluable resource for developing IoT intrusion detection systems.

### 3.2 Data Preprocessing

The most important part in the phase of model development. In this step the data is cleaned which means removing null values from the data and fixing the outliers and filling in missing values. The quality of data is improved for better classification and also improves the performance of the model. In the research the attack column indicates large number of unique values in the target which results in multiclass classification and to handle this large amounts of data, extra computational resources are required. The attack column is categorized based on the type of attack classes like DDoS, DoS, Reconnaissance, scanning, MITM and categories such as network attacks, exploitation attacks and some benign are categorized as unknown. And all the data is merged to form a final dataset for further evaluation. Data such as IP is converted into decimals and null values are checked.

The data preprocessing for CICIoT dataset was designed to make sure that high quality data which can support robust machine learning models for intrusion detection and other security applications. Provided the extensive volume of raw traffic data (548 GB) the preprocessing pipeline was essential in transforming the data into structured machine learning format. The large PCAP files was split into smaller 10 MB chunks for efficient processing. The DPKT library was used to extract 47 features from the network traffic. During the feature extraction the incomplete packets with null values were removed to make sure the data was clean and in usable format. The pre-processed data was stored in CSC files for easy access and analysis.

The raw PCAP files were retained for researchers to perform more complicated research if they intended to. The entire process was to make sure that the resulting dataset was well structured, balanced and suitable for use in various machine learning models which enables accurate and reliable detection and classification of IoT network traffic.

### **3.3 Exploratory Data Analysis**

The exploratory data analysis (EDA) plays a very major role as it helps in revealing the insights from the data. For UQ-NIDS, These insights are instrumental to visually examine various questions and patterns within the data. There are histogram to depict the distribution of IPV4\_SRC in relation to Attack Category. This also indicates high prevalence of exploitation attacks on IPV4\_SRC with lower chances of benign activities and also suggests that exploitation attacks on IPV4\_SRC should be prioritized for prevention or control.

For CICIoT dataset, it contains a total of 46,686,579 records including 33 distinct attack types categorized into seven major classes: DDoS, DoS, reconnaissance, web based attacks, brute force, spoofing and mirai botnet attacks. It also includes 1,098,195 records of benign traffic representing normal IoT device operations during both idle and interactive states. DDoS attacks dominate the dataset with over 33.9 million records. DoS attacks account for 8 million records while web-based and brute force attacks contribute 24,829 and 13,064 records respectively. The analysis of features reveals the distribution and importance of 47 extracted attributes such as flow duration, packet size, protocol type and flag counts. These features capture the unique characteristics of IoT network traffic which enables machine learning models to distinguish between benign and malicious activity. Correlation analysis is performed to identify relationships between features and reveal that certain attributes such as packet size and average packet length which are strongly correlated and may indicate redundancy.

The dataset migrates this through feature aggregation offering balanced representation for low-volume attack classes. Overall the EDA highlights the dataset's diversity, comprehensiveness for IoT security research. And by understanding the data's structure and attack patterns, researchers can effectively utilize the CICIoT dataset for developing and evaluating machine learning models to detect and classify IoT network threats. This analysis reinforces the dataset's value as realistic resource for advancing intrusion detection systems in IoT environments.

### **3.4 Feature Engineering**

Feature Engineering is very crucial for this research study, as it involves selecting features and transforming raw data into relevant information for machine learning models. And by using feature engineering the performance of algorithms including accuracy can be enhanced. And this process the data is first label-encoded, converting categorical features into numerical ones. Features such as IPV4\_SRC, IPV4\_DST, L4\_SRC, F4\_DST\_PORT and Attack Category are label encoded and the target column is subsequently dropped from this dataset and this is for UQ-NIDS dataset.

Whereas for the CICIoT dataset the feature engineering process begins with capturing network traffic with the help of Wireshark and storing it in PCAP format. This raw traffic is then processed using DPKT library to extract 47 critical features which are essential for understanding the behaviour of IoT devices under both benign and malicious scenarios. The

extracted features cover a broad range of the attributes including flow characteristics such as flow duration, packet-level details like packet size and protocol specific indicators. And the statistical summaries such as mean, standard deviation and variance are calculated for each flow along magnitude, radius and covariance which capture variability in packet lengths. These features are then aggregated into fixed size windows of 10 or 100 packets depending on the attack type and this approach helps to address data imbalance as high-volume attacks like DDoS generate significantly more traffic compared to web based or spoofing attacks. In order to further streamline the data the incomplete packets with null features are removed during preprocessing and makes sure that only clean and usable data is included.

Besides the timestamp data is excluded from feature sets as they do not contribute to network behaviour analysis. The CSV files provide structured and simplified dataset for machine learning which enables researchers to train models effectively. This feature engineering approach is built on best practices in IoT security and makes sure CICIoT dataset is flexible and supports pre-defined and custom feature extraction for advanced analysis The process also demonstrates a vigorously methodology for generating insights into IoT network behavior.

### **3.4.1 Feature Extraction**

This is the process of extracting the necessary features for machine learning. For the UQ-NIDS dataset, The features is important as they play an important role in the performance of machine learning which are extracted from the dataset. The dataset is split into training and testing set (train set is 80% & test size is 20%). And for extraction of the features Random Forest classifier is used and the important features are extracted. The top 8 features with highest scores are IPV4\_SRC and IPV4\_DST are then selected for the training of algorithm and these features are selected and normalized with the help of MinMaxScaler.

For the CICIoT dataset, The dataset is generated from simulating attacks on network of 105 IoT devices and during this simulation, network traffic is captured which records packets which are sent and received during both attack and benign scenarios. The data splitting procedure involved combining all captured malicious and benign traffic into single integrated dataset. This integrated dataset was then shuffled to make sure randomness and mitigate any inherent ordering from collection process.

The dataset was divided into 80% training data and 20% testing data. This split makes sure that training split provides sufficient data diversity for the machine learning model to learn pattern effectively while the testing set provides enough representation of unseen data to evaluate the generalization capabilities of the trained models.

All features were normalized using Standard Scaler method and this method ensures uniform range to improve performance and convergence of machine learning during training. The data split facilitated three classification tasks:

Binary classification: which distinguished between malicious and benign traffic.

Grouped classification: which categorized traffic into seven attack groups such as DDoS, DoS and reconnaissance.

Multiclass classification: which classified traffic into 33 individual attack types.

### 3.5 Model Training

This is the crucial part of making predictions and after identifying to extract 8 relevant features and normalized using minmax scaler. The dataset is divided into training and testing sets with 20% allocated for testing. The training and testing data are then reshaped into 2 dimensional formats. Four learning algorithms are used to classify the attacks. Logistic Regression, XG Boost, Random Forest and Decision Tree. The training data is then processed with the above algorithms.

For the UQ-NIDS dataset, Decision Tree classifier using ‘scikit-learn’ library for learning task and evaluate performance through metrics. The libraries are imported for building the model. Decision Tree works by recursively partitioning the feature space into smaller subsets based on feature values which will form a tree like structure where each node represents decision rule on feature and each leaf node corresponds to predicted class. For the UQ-NIDS dataset, Decision Tree is trained on the reshaped training data (X\_train\_reshaped) and corresponding class labels (Y\_train\_original ). The next model we used is Random Forest, an this algorithm builds multiple decision trees during training and combines their predictions improve accuracy and robustness.

In our program, Random Forest operates by creating multiple decision trees, each trained on random subset of training data, sampled with replacement and for each split in a tree the algorithm considers random subset of the features rather than all features which includes further randomness. This dual randomness ensures diversity among the trees. The next model which we have implemented to UQ-NIDS is XGBoost, as this is highly efficient and scalable machine learning model based on gradient boosting framework. XGBoost works by building on decision trees where each subsequent tree aims to correct the errors made by the previous trees. The last model implemented to UQ-NIDS is Logistic Regression, which is simple but powerful linear classification algorithm to classify input data. This algorithm models the relationship between input features and target classes by fitting a logistic function to the weighted sum of the input features. As this ensures that output probabilities lie between 0 and 1 which are then used to classify data into distinct categories.

For the CICIoT dataset, supports multi-level classification of IoT network traffic categorized into 34 class, 8 class and 2 class. The dataset includes 33 distinct attack types and one class for benign traffic. 34 class is the hardest classification task because of the high diversity in traffic patterns between different attack types. Also some of the classes have imbalanced data distributions as certain attacks generate much more traffic than others which leads to challenges. The 8 class classification the attacks are grouped into eight broader categories. This classification aggregates similar types of attacks into broader categories and simplifies the task. This classification is easier than the 34 class classification this still requires models to distinguish between overlapping traffic features in some categories. The 2 class classification is simply classified as Benign and Malicious. All attack types are grouped under malicious class. This is the simplest task because it only requires distinguishing traffic from normal traffic.

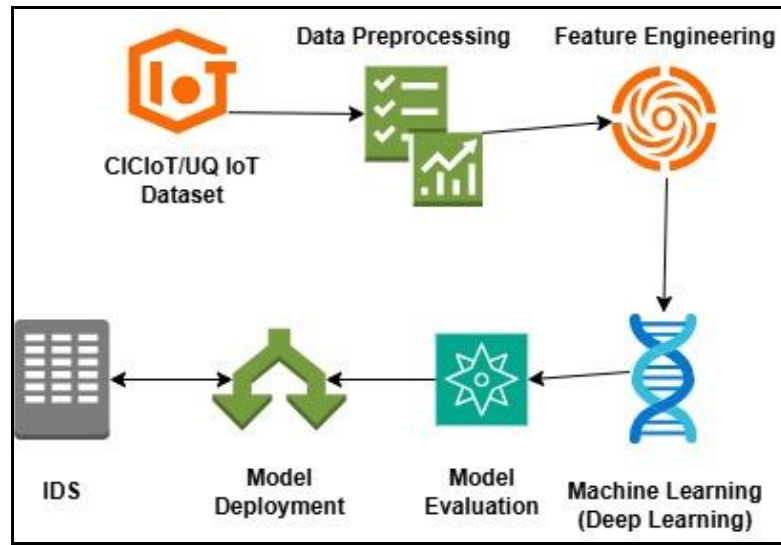
## 4 Design Specification

As we have previously discussed data collection, preprocessing, feature engineering, model training in the methodology section. Our IDS makes use of deep learning model trained to classify traffic packets as either normal or potential malicious and in case of CICIoT, it is

divided into three different classifications. When the IDS classifies network packets, it triggers alert if malicious activity is detected detailing the attack so that admin and user can respond. Normal traffic is allowed to pass through which makes sure server can continue providing service to genuine users.

The UQ-NIDS has been trained for Decision Tree, Random Forest, XGBoost, Logistic Regression and Naïve Bayes. And we have implemented hyper parameter tuning in order to select which is the optimal set for the given model. The model improved its performance as a result of proper tuning of hyperparameters. And hyperparameters helps in finding the right balance between overfitting and underfitting.

The framework presented in this paper illustrates how the system would hypothetically operate as an IDS in aim to provide insights to potential deployment and operation of the system. For the CICIoT dataset is optimised for machine learning evaluations by extracting significant traffic features and maintaining flexibility for feature engineering. The framework which demonstrates the functionalities of IDS would operate and it offers insights to its operation. The CICIoT dataset is well optimized for its machine learning evaluations as it extracts traffic features and maintains its flexibility for feature engineering and this in turn allows IDS to adapt and improve the detection capabilities. The CICIoT dataset is classified into three types based on the attack categories. 34 class (33+1) classification (benign and all individual attacks), 8 class (7+1) classification (benign and attack categories) and 2 class classification (1+1) classification (benign and malicious).



## 5 Evaluation

Our research project utilized Python which is a versatile programming language particularly suited for data science and machine learning tasks. We have employed several Python libraries through out the project and numerical computations as we have used NumPy to handle large multi-dimensional arrays and perform important mathematical operations. We have pandas for data preprocessing, cleaning and structuring both UQ-NIDS and CICIoT datasets into format suitable for analysis. In order to visualize the data trends and assess the models' performance we have made use of Plotly for creating interactive plots and Matplotlib

for static visualizations. Scikit-learn (sklearn) was used to implement machine learning models which offers a wide range of algorithms.

Bambang Susilo et al. (2020) suggested that Random Forest achieves high AUC values for multiclass classification and best suited as it combines decision trees to classify data and offers higher accuracy for multiclass classification.

Our research has also performed hyperparameter tuning to optimize the performance and capabilities. When the model becomes very complex and memorizing the training data and failing to implement it on new data. Parameters like the number of layers and maximum depth in decision trees control the model complexity.

Without Hyperparameter tuning:

| Model (For <b>UQ-NIDS</b> ) | Accuracy | Precision | Recall | F1 Score |
|-----------------------------|----------|-----------|--------|----------|
| Decision Tree               | 96.24    | 97.15     | 97.13  | 97.13    |
| Random Forest               | 97.12    | 97.15     | 97.12  | 97.12    |
| XG Boost                    | 97.08    | 97.18     | 97.26  | 97.09    |
| LogReg                      | 64.59    | 65.65     | 64.59  | 64.67    |
| Naive Bayes                 | 49.47    | 56.37     | 49.47  | 46.00    |

With Hyperparameter tuning:

| Model (With Hyper Parameter Tuning) | Accuracy | Precision | Recall | F1 Score |
|-------------------------------------|----------|-----------|--------|----------|
| Decision Tree                       | 96.74    | 96.74     | 96.74  | 96.74    |
| Random Forest                       | 97.04    | 97.04     | 97.12  | 97.12    |
| XG Boost                            | 97.08    | 97.18     | 97.08  | 97.09    |
| LogReg                              | 68.87    | 68.96     | 68.87  | 68.3     |
| Naive Bayes                         | 55.95    | 61.96     | 55.95  | 55.68    |

The models performed better with hyper parameter tuning. Without the HP tuning, Decision Tree achieved 96.24, 97.15, 97.13, 97.13 and with HP tuning 96.74 throughout. And as for Logistic Regression without HP tuning, it achieved 64.59, 65.65, 64.59 and 64.67 and with HP tuning it received better results 68.87, 68.96, 68.87 and 68.3. Overall the performances improved when the models were subjected to hyper parameter tuning and these results helped with our research to attain better results.

The evaluation process for the CICIoT dataset involves structured pipeline to assess the performance of various machine learning methods for detecting and classifying the malicious IoT network traffic. The main focus is on the three levels of classification: binary (benign and malicious), grouped (7 attack category), multiclass (34 individual attacks). The dataset is initially combined, shuffled and divided into training and testing subsets. The features are then normalized using standard scaler to make sure consistent scaling. The evaluation highlighted the robustness.

In cases like Logistic Regression on the UQ-NIDS dataset, where accuracy (68.87%) had conflicted with low precision (68.96%) and recall (68.87%), thus the focus on F1 score helped balance these metrics. Introducing the hyperparameter tuning were to include adjustments to learning rates and decision thresholds further aligned with performance for metrics. And for models such like the support vector machines which had exhibited a high accuracy but low recall in multiclass classification the class weights were adjusted to prioritize minority classes which enhanced overall performance.

On the CICIoT dataset the feasibility of multiclass classification at 34-class, 8- class, and 2-class levels was established, with Decision Tree consistently delivering the best results across tasks. This research also provided impactful insights and confirming that hyperparameter tuning and tailored the dataset specific preprocessing significantly to enhance model performance. And ensemble methods like Random Forest and XGBoost proved their robust for IoT anomaly detection.

For 34 (33+1) class classification:

| Model (For <b>CICCIoT</b> ) | Accuracy | Precision | Recall  | F1 Score |
|-----------------------------|----------|-----------|---------|----------|
| Decision Tree               | 99.1981  | 80.6247   | 81.0326 | 80.6382  |
| Random Forest               | 99.1643  | 70.4492   | 83.1586 | 71.40.21 |
| LogReg                      | 80.1831  | 59.4978   | 48.528  | 49.1284  |
| Support Vector Machine      | 78.7127  | 52.8674   | 42.7667 | 43.3722  |

For 8 (7+1) class classification:

| Model (For <b>CICCIoT</b> ) | Accuracy | Precision | Recall  | F1 Score |
|-----------------------------|----------|-----------|---------|----------|
| Decision Tree               | 99.4054  | 83.1231   | 82.757  | 82.9336  |
| Random Forest               | 99.4368  | 70.54     | 91      | 71.9289  |
| LogReg                      | 83.147   | 51.0881   | 68.3171 | 53.7237  |
| Support Vector Machine      | 82.3073  | 67.7474   | 46.8986 | 50.1793  |

For 2 (1+1) class classification:

| Model (For <b>CICCIoT</b> ) | Accuracy | Precision | Recall  | F1 Score |
|-----------------------------|----------|-----------|---------|----------|
| Decision Tree               | 99.5888  | 95.5463   | 95.5038 | 95.525   |
| Random Forest               | 99.68    | 96.5395   | 96.5163 | 96.5279  |
| LogReg                      | 98.902   | 86.3226   | 89.0443 | 87.6315  |
| Support Vector Machine      | 98.7115  | 87.1835   | 83.5327 | 85.2584  |

## 5.1 Evaluation based on Accuracy

Accuracy is used to display the proportion of instances that are classified correctly out of the total instances. This is used when all the classes are balanced. The accuracy and model is

directly proportional to each other, in simpler terms higher the accuracy and the model is said to perform better. In our project, the evaluation of each model of accuracy provides deeper insight to the performance of the model. Comparative analysis of models is given in table.

For UQ-NIDS dataset, we have implemented hyperparameter tuning in order to improve the model efficiency. The learning rate is adjusted so it can affect how quickly the model can coverage during the training process. Hyperparameters tuning also aims to optimize the performance of the model. From the above table it is shown that XGBoost (97.26%) has the highest accuracy and Naïve Bayes (55.71%) has the lowest accuracy. Hyperparameter tuning improves the model performance by adjusting the learning rate as in order to enhance the efficiency and training of the model. XGBoost showed strong performance on the UQ-NIDS dataset while the Naïve Bayes underperformed with comparison to it.

The CICIoT dataset took about 12 hours to completely execute the program as they more that 12 GB of storage which the system had to handle. And for the CICIoT dataset, the classification plays a vital role. And for 32 class classification (33+1), the accuracy is highest for Decision Tree (99.19%) and Support Vector Machine (78.17%) has the lowest accuracy.

For the 8 class classification of the CICIoT dataset, the highest accuracy is achieved by Decision Tree (99.40%) and the lowest accuracy for the 8 class classification belongs to Support Vector Machine model (82.30%).

The 2 (1+1) class classification of the CICIoT dataset. Decision Tree (99.58%) holds the highest accuracy and Support Vector Machine (98.71%) has the lowest accuracy. Decision Tree consistently outperformed other models in various classification tasks and achieves the highest accuracy across 32, 8 and 2 class classification. Support Vector Machine generally had lower accuracy which indicates it may not be the suitable model for this classification tasks.

## **5.2 Evaluation based on Precision**

This is used to tell the proportion of positive prediction which are correct and it is used when the rate of false positive is high. In our experiment for the UQ-NIDS dataset, XGBoost gave the highest precision (97.36%) and the lowest was achieved by Naïve Bayes (62.15%). With hyperparameter tuning, the models performed XGBoost (97.08%).

For CICIoT dataset, the classification is divided into 3 categories. 34 (33+1) class classification based on the attack categories. The highest precision is achieved by Random Forest (96.55%) and the lowest achieved model is Support Vector Machine (52.86%).

For the 8 class classification of the CICIoT dataset, the highest precision is achieved by Decision Tree (99.40%) and the lowest precision for the 8 class classification belongs to Support Vector Machine model (82.30%).

For the 2 (1+1) class classification of the CICIoT dataset, the highest precision is achieved by Decision Tree (99.40%) and the lowest precision for the 8 class classification belongs to Support Vector Machine model (82.30%).



### 5.3 Evaluation based on Recall

Recall also known as True Positive rate tells the proportion of actual positive that are correctly identified, and recall is used when the rate of false positive is high. In our experiment the recall metric is used to get the insights of models' ability to get the instance which are positive.

The UQ-NIDS dataset, XGBoost gave the highest recall (97.36%) and the least recall was achieved by Naïve Bayes (62.15%). With hyperparameter tuning, the models performed as the same with little improvements.

For CICIoT dataset, the classification is divided into 3 categories. 34 (33+1) class classification based on the attack categories. The highest recall is achieved by Random Forest (96.44%) and the lowest recall achieved model is Support Vector Machine (42%).

For the 8-class classification of the CICIoT dataset, the highest precision is achieved by Decision Tree (99.40%) and the lowest precision for the 8 class classification belongs to Support Vector Machine model (82.30%).

For the 2 (1+1) class classification of the CICIoT dataset, the highest precision is achieved by Decision Tree (99.40%) and the lowest precision for the 8 class classification belongs to Support Vector Machine model (82.30%).

### 5.4 Evaluation based on F1-Score

In our experiment, F1 score tells the harmonic mean of precision and recall as it provides a single metric that balances both the measures. This study also focuses on selecting features crucial for evaluation and model training. This performance demonstrates its effectiveness in maintaining a high balance between precision and recall, making it as the most reliable among others.

The UQ-NIDS dataset, XGBoost gave the highest f1-score (97.36%) and the lowest f1-score was achieved by Naïve Bayes (62.15%). With hyperparameter tuning, the models performed with little or not much difference result metrics.

For CICIoT dataset, the classification is divided into 3 categories. 34 (33+1) class classification based on the attack categories. The highest recall is achieved by Random Forest (96.44%) and the lowest recall achieved model is Support Vector Machine (43.3722%).

For the 8-class classification of the CICIoT dataset, the highest precision is achieved by Decision Tree (99.40%) and the lowest precision for the 8 class classification belongs to Support Vector Machine model (82.30%).

For the 2 (1+1) class classification of the CICIoT dataset, the highest precision is achieved by Decision Tree (99.40%) and the lowest precision for the 8-class classification belongs to Support Vector Machine model (82.30%).

## 5.5 Discussion

In this research five different algorithms for intrusion detection system (IDS) have been implemented to two different datasets in order to analysis in depth which model is more suitable. The algorithms for CICIoT dataset took more than 12 hours to execute as there were three different types of attack classification and the one of the major objective of our research was to figure out which machine learning and deep learning model is best suited for these attack categories namely the 34 (33+1) class, 8 (7+1) class and the binary classification which is the 2 (1+1) class classification and provide the AUC results to compare which model produced the best results.

The models used are Logistic Regression, Support Vector Machine, Decision Tree, Random Forest and XGBoost for classifying these attacks. And the performance of these algorithms are evaluated using the four key metrics: accuracy, precision, recall and F1 score to provide comparison. Among the algorithms trained on UQ-NIDS dataset, XGBoost model outperformed Random Forest, Decision Tree and Logistic Regression. The results are almost the same after implementing hyperparameter tuning.

The Decision Tree , models achieved higher accuracies respectively. The research results highlights the importance of feature selection for training the model, indicating Decision Tree and Logistic Regression is a reliable choice for 34 class.

For 8 class classification Decision Tree and Random Forest are the two best models followed by Logistic Regression and Support Vector Machine.

For 2 class classification, once again Decision Tree outperformed the other models. This emphasizes the significance of considering machine learning models for complex classification tasks as they improve the overall performance over traditional approaches.

Our research compared to previous work only for the CICIoT dataset is shown in the below table. The machine learning models Decision Tree, Random Forest, Naïve Bayes and Logistic Regression along with hyperparameter tuning for each were not implemented on UQ-NIDS for their previous work.

| Model            | Accuracy | Recall  | Precision | F1-Score |
|------------------|----------|---------|-----------|----------|
| 34 class Log Reg | 80.2351  | 59.5201 | 48.6752   | 49.3884  |
| 8 class Log Reg  | 83.1674  | 69.6055 | 51.2409   | 53.9424  |
| 2 class Log Reg  | 98.9023  | 89.04   | 86.3157   | 87.6258  |

*Table: Previous past paper*

| Model            | Accuracy | Recall  | Precision | F1-Score |
|------------------|----------|---------|-----------|----------|
| 34 class Log Reg | 80.1831  | 59.4978 | 48.528    | 49.124   |
| 8 class Log Reg  | 83.147   | 68.3171 | 51.0881   | 53.7237  |
| 2 class Log Reg  | 98.902   | 89.0443 | 86.3226   | 87.6315  |

*Comparison Table: Our Research Paper*

The above table depicts the comparison between our research paper and previous research paper. From the previous paper, it is evident that only Logistic Regression is being carried out and the rest of the models are novel contributions to CICIoT dataset.

## 6. Conclusion and Future Work

There are plenty of promising avenues for the future work. The hybrid model could be further developed and implemented to enhance the efficiency and to handle such large scale and data streams much more efficiently. And incorporating new online learning would enable the model to be in continually adapt new cyber threats. These enhancements can lead to conduct extensive performance evaluations for in various real-world environments and IoT platforms as they would provide more evidence of the model's suitability across the different deployments and this can promote broader use and can be potentially advance the cybersecurity defenses further.

Features were selected based on the contribution of them differentiating between benign and malicious traffic. The major key attributes such as the flow duration, packet size and protocol indicators were also prioritized for their direct relevance to the anomalous network behaviors. And the statistical attributes like mean, variance and as well as standard deviation were also included to provide comprehensive understanding of this data. And features which were unrelated to network behavior like timestamps were removed and introduced as noise. And highly correlated features are also removed through correlation analysis which can reduce redundancy. Computationally intensive features that has minimal impact on classification accuracy have been taken out in order to enhance the efficiency. The prioritization criteria were more focused on the relevance and selecting features like the flow duration and packet size for their ability to capture traffic.

## Reference

Guezzaz, A., Benkirane, S., & Azrou, M. (2022). A Novel Anomaly Network Intrusion Detection System for Internet of Things Security. *EAI/Springer Innovations in Communication and Computing*, 129–138. [https://doi.org/10.1007/978-3-030-90083-0\\_10](https://doi.org/10.1007/978-3-030-90083-0_10)

Kayode Saheed, Y., Idris Abiodun, A., Misra, S., Kristiansen Holone, M., & Colomo-Palacios, R. (2022). A machine learning-based intrusion detection for detecting internet of things network attacks. *Alexandria Engineering Journal*, 61(12), 9395–9409. <https://doi.org/10.1016/J.AEJ.2022.02.063>

Vykopal, J.; Plesnik, T.; Minarik, P. Network-based dictionary attack detection. In *Proceedings of the 2009 International Conference on Future Networks*, Bangkok, Thailand, 7–9 March 2009; pp. 23–27.

Lata, S., & Singh, D. (2022). Intrusion detection system in cloud environment: Literature survey & future research directions. *International Journal of Information Management Data Insights*, 2(2), 100134. <https://doi.org/10.1016/J.IJIMEI.2022.100134>

Layeghy, S., Gallagher, M., & Portmann, M. (2024). Benchmarking the benchmark — Comparing synthetic and real-world Network IDS datasets. *Journal of Information Security and Applications*, 80, 103689. <https://doi.org/10.1016/J.JISA.2023.103689>

Neto, E.C.P.; Dadkhah, S.; Ghorbani, A.A. Collaborative DDoS Detection in Distributed Multi-Tenant IoT using Federated Learning. In *Proceedings of the 2022 19th Annual*

International Conference on Privacy, Security & Trust (PST), Fredericton, NB, Canada, 22–24 August 2022; pp. 1–10.

Kang, H.; Ahn, D.H.; Lee, G.M.; Yoo, J.; Park, K.H.; Kim, H.K. IoT network intrusion dataset. *IEEE Dataport* 2019,1, 1.

Long, Z., Yan, H., Shen, G., Zhang, X., He, H., & Cheng, L. (2024). A Transformer-based network intrusion detection approach for cloud security. *Journal of Cloud Computing*, 13(1), 1–11. <https://doi.org/10.1186/S13677-023-00574-9/TABLES/4>

Safi, M.; Dadkhah, S.; Shoeleh, F.; Mahdikhani, H.; Molyneaux, H.; Ghorbani, A.A. A Survey on IoT Profiling, Fingerprinting, and Identification. *ACM Trans. Internet Things* 2022, 3, 1–39.

Nizamudeen, S. M. T. (2023). Intelligent intrusion detection framework for multi-clouds – IoT environment using swarm-based deep learning classifier. *Journal of Cloud Computing*, 12(1), 1–14. <https://doi.org/10.1186/S13677-023-00509-4/FIGURES/6>

Kaur, B.; Dadkhah, S.; Xiong, P.; Iqbal, S.; Ray, S.; Ghorbani, A.A. Verification based scheme to restrict iot attacks. In *Proceedings of the 2021 IEEE/ACM 8th International Conference on Big Data Computing, Applications and Technologies (BDCAT'21)*, Leicester, UK, 6–9 December 2021; pp. 63–68.

Ferrag, M.A.; Friha, O.; Hamouda, D.; Maglaras, L.; Janicke, H. Edge-IIoTset: A new comprehensive realistic cyber security dataset of IoT and IIoT applications for centralized and federated learning. *IEEE Access* 2022, 10, 40281–40306

Pate, J.; Adegbiya, T. AMELIA: An application of the Internet of Things for aviation safety. In *Proceedings of the 2018 15th IEEE Annual Consumer Communications & Networking Conference (CCNC)*, Las Vegas, NV, USA, 12–15 January 2018; pp. 1–6.

Zantalis, F.; Koulouras, G.; Karabetsos, S.; Kandris, D. A review of machine learning and IoT in smart transportation. *Future Internet* 2019, 11, 94.

deRito, C.; Bhatia, S. Comparative Analysis of Open-Source Vulnerability Scanners for IoT Devices. In *Intelligent Data Communication Technologies and Internet of Things*; Springer: Singapore, 2022; pp. 785–800.