



# Cyber Threat Detection using Machine Learning

**Nibedita Baidya**

Department of Computer Science and Engineering (Artificial Intelligence)  
GIFT Autonomous, Bhubaneswar, Odisha, India

**Swarnila Mallick**

Department of Computer Science and Engineering GIFT Autonomous, Bhubaneswar, Odisha, India

**Debaprasad Nanda**

Department of Computer Science and Engineering  
GIFT Autonomous, Bhubaneswar, Odisha, India

## ABSTRACT

Cybersecurity threats are increasing rapidly with the growth of internet-based services, cloud computing, and digital communication systems. Traditional security methods are often unable to detect sophisticated and evolving cyber-attacks in real time. To overcome these limitations, this project presents a Cyber Threat Detection System using Machine Learning techniques for identifying malicious activities and improving network security. The system is designed to analyze large volumes of network traffic and detect abnormal behavior that may indicate cyber threats such as phishing attacks, malware, denial-of-service attacks, unauthorized access, and data breaches.

The proposed system uses machine learning algorithms to classify network activities into normal and malicious categories based on patterns extracted from the dataset. Data preprocessing techniques such as data cleaning, feature selection, normalization, and encoding are applied to improve the quality and accuracy of the model. Various machine learning algorithms including Decision Tree, Random Forest, Logistic Regression, and Support Vector Machine are used for threat prediction and comparison of performance. The system is trained and tested using cybersecurity datasets to achieve accurate detection results with reduced false positives.

In addition to threat detection, the project integrates data visualization techniques using Tableau dashboards to provide real-time graphical analysis of cyber threats, attack trends, and network behavior. The dashboard helps administrators monitor suspicious activities effectively and make faster security decisions. Performance evaluation metrics such as accuracy, precision, recall, and F1-score are used to measure the effectiveness of the proposed model.

## 1. INTRODUCTION

This project, “Cyber Threat Detection Using Machine Learning,” focuses on developing an intelligent system capable of detecting malicious activities within network traffic using machine learning algorithms. The system analyzes large cybersecurity datasets to identify suspicious patterns and classify network behavior as either normal or malicious. Different machine learning algorithms such as Decision Tree, Random Forest, Logistic Regression, and Support Vector Machine are implemented and evaluated to determine the most effective model for cyber threat detection.

The project also emphasizes the importance of data preprocessing techniques, including data cleaning, feature selection, normalization, and data transformation, to improve model performance and prediction accuracy. By training the model on historical network traffic data, the system can recognize abnormal activities and generate alerts for potential cyber-attacks in real time.

In addition to threat detection, the project integrates Tableau dashboards for data visualization and analysis. The dashboard provides graphical representations of cyber threats, attack frequency, traffic behavior, and detection results, helping administrators monitor network activities more effectively. Visualization tools improve decision-making by presenting complex cybersecurity data in an understandable and interactive format.

The main aim of this project is to build a reliable, scalable, and efficient cyber threat detection system that can reduce manual monitoring efforts and enhance network security. The implementation of machine learning and visualization technologies demonstrates how artificial intelligence can contribute to modern cybersecurity solutions by improving detection speed, accuracy, and overall system performance. This project provides practical knowledge in machine learning, cybersecurity, data analytics, and visualization, making it highly relevant in today’s digital era



where protecting sensitive information has become a critical requirement.

The project also includes data preprocessing and feature extraction techniques to improve the quality of the dataset and enhance machine learning model performance. Important features such as IP addresses, protocol types, packet sizes, and traffic duration are analyzed to detect abnormal network behavior.

## 2. OBJECTIVE

The main objective of the project “Cyber Threat Detection Using Machine Learning and Tableau” is to develop an intelligent and automated cybersecurity system capable of detecting malicious activities and cyber threats in network traffic using machine learning algorithms and data visualization techniques.

The project also aims to implement advanced machine learning algorithms such as Neural Network and Support Vector Classifier (SVC) for cyber threat detection. These algorithms are trained using historical network traffic data to identify cyber-attacks such as phishing attacks, malware infections, denial-of-service attacks, intrusion attempts, and unauthorized access. The objective is to achieve high detection accuracy, reduce false alarms, and improve the efficiency of cybersecurity monitoring.

Another major objective of the project is to automate the cyber threat detection process and reduce manual intervention in network security management. The system is designed to generate alerts and notifications whenever suspicious activities are detected so that administrators can respond quickly to security incidents. This helps improve incident response time and minimizes the risk of data breaches and system damage.

1. To visualize cyber threat data using Tableau dashboards
2. To implement machine learning algorithms for threat detection
3. To improve the accuracy and efficiency of cyber threat detection
4. To automate the cybersecurity monitoring process
5. To perform data preprocessing and feature engineering
6. To compare the performance of different machine learning models

## 3. LITERATURE SURVEY

Cybersecurity has become one of the most critical areas of research due to the rapid increase in cyber-attacks, data breaches, phishing, malware, ransomware, and unauthorized network access. Traditional security systems such as firewalls and signature-based intrusion detection systems are often unable to detect new and sophisticated attacks effectively. As a result, researchers have focused on developing intelligent cyber threat detection systems using Machine Learning (ML) and Artificial Intelligence (AI) techniques. Several studies have been conducted to improve intrusion detection accuracy, reduce false alarms, and identify unknown threats in real-time environments.

Research published in Springer Journal of Big Data discussed machine learning-based intrusion detection systems and compared popular cybersecurity datasets such as NSL-KDD, CICIDS2017, and UNSW-NB15. The study identified challenges related to dataset imbalance, outdated attack patterns, and model generalization. The authors emphasized the importance of modern datasets and real-time network analysis for building effective cyber threat detection systems.

Another survey focused on the CSE-CIC-IDS2018 dataset and analyzed various intrusion detection models developed using machine learning techniques. The research explained that anomaly-based detection systems perform better than traditional signature-based methods because they can identify previously unknown attacks. The paper also highlighted the role of big data analytics in cybersecurity and the increasing use of advanced ML algorithms for network traffic analysis.

A survey on malicious URL detection using machine learning explained how cybercriminals use malicious websites and phishing links to attack users. The researchers proposed machine learning-based detection models capable of identifying harmful URLs using feature extraction and classification techniques. The study showed that ML models can detect newly generated malicious websites more effectively than traditional blacklist-based systems.

## 4. EXISTING SYSTEM

The existing system for cyber threat detection mainly depends on traditional security mechanisms such as firewalls, antivirus software, signature-based intrusion detection systems (IDS), and Security Information and Event Management



(SIEM) tools. These systems are designed to detect known threats by comparing network activities and files against predefined rules, attack signatures, or malware databases. Traditional cybersecurity solutions are widely used in organizations to monitor network traffic, block unauthorized access, and protect sensitive information from cyber-attacks.

In the existing approach, signature-based detection is one of the most commonly used methods. This method identifies threats only when the attack pattern already exists in the database. Whenever a new malware signature or attack type is discovered, the security database must be updated manually or through regular updates. Although this approach is effective for detecting previously known threats, it fails to identify zero-day attacks, polymorphic malware, and evolving cyber threats that continuously change their behavior to avoid detection.

Another existing approach is rule-based intrusion detection systems, where predefined rules are created to monitor suspicious network activities. These systems generate alerts whenever unusual behavior matches the configured rules. However, rule-based systems produce a large number of false positives and false negatives, making it difficult for security analysts to identify real threats accurately. Additionally, these systems require continuous monitoring and manual rule updates to remain effective against modern cyber-attacks.

Traditional SIEM systems are also widely used for collecting and analyzing security logs from multiple devices and applications. SIEM platforms help organizations centralize security monitoring and incident reporting. However, conventional SIEM systems often struggle with scalability, real-time threat analysis, alert fatigue, and lack of intelligent automation. Modern attacks involving AI-powered phishing, fileless malware, and advanced persistent threats can bypass these systems because they rely heavily on static rules and historical attack patterns.

Overall, the existing cyber threat detection systems suffer from several limitations such as dependence on known attack signatures, inability to detect unknown threats, high false alarm rates, lack of automation, poor scalability, delayed response time, and limited real-time analysis capabilities. These drawbacks highlight the need for intelligent cybersecurity solutions using machine learning and data visualization techniques that can detect evolving threats more accurately and efficiently.

## 5. PROPOSED SYSTEM

The existing system for cyber threat detection mainly depends on traditional security mechanisms such as firewalls, antivirus software, signature-based intrusion detection systems (IDS), and Security Information and Event Management (SIEM) tools. These systems are designed to detect known threats by comparing network activities and files against predefined rules, attack signatures, or malware databases. Traditional cybersecurity solutions are widely used in organizations to monitor network traffic, block unauthorized access, and protect sensitive information from cyber-attacks.

In the existing approach, signature-based detection is one of the most commonly used methods. This method identifies threats only when the attack pattern already exists in the database. Whenever a new malware signature or attack type is discovered, the security database must be updated manually or through regular updates.

Another existing approach is rule-based intrusion detection systems, where predefined rules are created to monitor suspicious network activities. These systems generate alerts whenever unusual behavior matches the configured rules.

Traditional SIEM systems are also widely used for collecting and analyzing security logs from multiple devices and applications. SIEM platforms help organizations centralize security monitoring and incident reporting. However, conventional SIEM systems often struggle with scalability, real-time threat analysis, alert fatigue, and lack of intelligent automation. Modern attacks involving AI-powered phishing, fileless malware, and advanced persistent threats can bypass these systems because they rely heavily on static rules and historical attack patterns.

## 6. SYSTEM REQUIREMENT

The hardware requirements of the project include a computer system with a minimum Intel Core i3 processor or equivalent AMD processor to handle machine learning computations and data processing tasks. A minimum of 4 GB RAM is required for running Python libraries, preprocessing datasets, and performing machine learning operations, although 8 GB RAM or higher is recommended for better performance and faster processing of large cybersecurity datasets. The system should have at least 250 GB of storage space for storing datasets, trained models, project files, reports, and visualization dashboards. A stable internet connection is also required for downloading datasets, installing software



packages, and accessing online resources related to cybersecurity and machine learning.

- Software Requirements
- Operating System : Windows 10/11, Linux, or macOS
- Database: MySQL / SQLite
- Dataset Format: CSV / Excel Files
- Python Version: Python 3.8 or higher
- Browser: Google Chrome / Microsoft Edge

## 7. SYSTEM ARCHITECTURE

The architecture begins with the data collection layer, where cybersecurity-related data is gathered from different sources such as network traffic, system logs, login records, IP addresses, protocols, and attack datasets. This collected data may contain raw and unstructured information that cannot be directly used for machine learning analysis. Therefore, the next stage of the architecture focuses on data preprocessing. In this stage, the raw data is cleaned by removing duplicate records, handling missing values, converting inconsistent formats, and normalizing the data to improve its quality and reliability. Proper preprocessing is very important because the performance of machine learning models depends heavily on the quality of the input dataset.

After preprocessing, the system performs feature extraction and feature selection. In this phase, important attributes such as packet size, protocol type, source IP address, destination IP address, traffic duration, and connection status are selected from the dataset. Selecting only relevant features helps reduce unnecessary data complexity and improves the efficiency and accuracy of the machine learning algorithms. The processed and optimized dataset is then sent to the machine learning layer, which acts as the core component of the system.

The machine learning layer is responsible for analyzing network behavior and identifying malicious activities. In this project, various machine learning algorithms such as Decision Tree, Random Forest, Logistic Regression, and Support Vector Machine are used to train and test the model. These algorithms learn patterns from historical cybersecurity data and classify network activities as either normal or malicious. The trained system can detect different types of cyber threats including phishing attacks, malware infections, denial-of-service attacks, intrusion attempts, and unauthorized access.

The final layer of the architecture is the visualization dashboard layer, where Tableau is used to create interactive dashboards and graphical reports. The dashboard displays important cybersecurity information such as attack statistics, threat categories, traffic analysis, suspicious activities, and model performance through charts, graphs, and reports. These visualizations make it easier for administrators and security analysts to understand complex cybersecurity data and make faster security decisions.

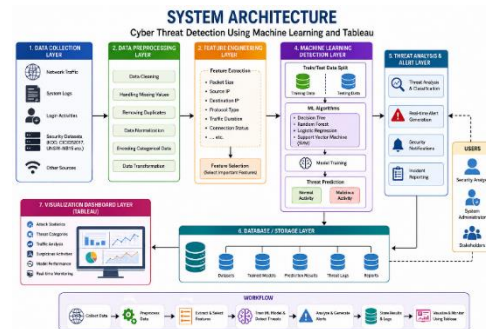


Fig 1: System Architecture Diagram

## 8. DATA FLOW DIAGRAM

The Data Flow Diagram (DFD) of the “Cyber Threat Detection Using Machine Learning” project represents the movement of data throughout the system and explains how information is processed from input to output. The DFD helps in understanding the flow of cybersecurity data between different modules, processes, databases, and users. It illustrates how network traffic data is collected, analyzed using machine learning algorithms, stored in databases, and visualized through dashboards for cyber threat monitoring and detection.

In the proposed system, the data flow begins with the data source layer, where network traffic data, system logs, IP addresses, login records, and cybersecurity datasets are collected from different sources. These datasets may include information related to normal and malicious activities such as phishing attacks, malware behavior, denial-of-service attacks, and intrusion attempts. The collected data acts as the input for the entire cyber threat detection system.

After collecting the data, it is transferred to the data preprocessing module. In this module, the raw cybersecurity data is cleaned and transformed into a structured format suitable for machine learning analysis. The preprocessing stage removes duplicate records, handles missing values, normalizes numerical values, and converts categorical data into machine-readable formats. This step is very important because high-quality



data improves the performance and accuracy of the machine learning models.

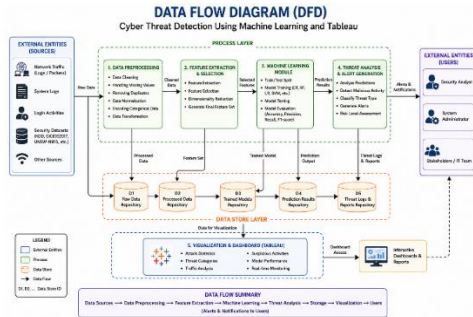


Fig 2: Data Flow Diagram

## 9. DATABASE DESIGN

The database design for the “Cyber Threat Detection Using Machine Learning and Tableau” project is developed to store, manage, and organize cybersecurity data efficiently. The database plays an important role in maintaining network traffic records, user information, machine learning prediction results, threat logs, and dashboard reports. A well-structured database design helps improve system performance, data retrieval speed, scalability, and overall threat analysis accuracy.

The Feature Extraction Table stores selected features extracted from the preprocessed data. This table includes fields such as Feature ID, Traffic ID, Feature Name, Feature Value, and Extraction Date. These features are later used by machine learning algorithms for threat detection and classification.

The Machine Learning Prediction Table stores the prediction results generated by the machine learning models. This table contains fields such as Prediction ID, Traffic ID, Algorithm Name, Prediction Result, Threat Type, Accuracy Score, and Prediction Time. The prediction result identifies whether the network activity is normal or malicious.

The Threat Log Table is responsible for storing detailed records of detected cyber threats and suspicious activities. This table contains fields such as Threat ID, Prediction ID, Threat Category, Risk Level, Detection Time, Alert Status, and Incident Description. The stored threat logs help in future analysis, reporting, and security investigations.

## 10. Module Description

### 10.1. Data Collection Module

The Data Collection Module is the first and one of the most important modules of the system. This module is responsible for gathering cybersecurity-

related data from various sources such as network traffic, system logs, login records, IP addresses, and publicly available cybersecurity datasets.

### 10.2. Data Preprocessing Module

The Data Preprocessing Module is responsible for cleaning and transforming the raw data into a structured format suitable for machine learning analysis. The collected data may contain missing values, duplicate records, inconsistent formats, and irrelevant information that can reduce model performance. Therefore, preprocessing is performed to improve data quality and reliability.

### 10.3. Feature Extraction and Selection Module

The Feature Extraction and Selection Module identifies the most important attributes from the processed dataset that are useful for detecting cyber threats. In cybersecurity datasets, many attributes may be irrelevant or redundant, which can increase system complexity and reduce prediction accuracy.

### 10.4. Threat Detection and Analysis Module

The Threat Detection and Analysis Module analyzes the prediction results generated by the machine learning models. If suspicious or malicious activity is detected, the system classifies the threat according to its type and severity level.

### 10.5. Alert and Notification Module

The Alert and Notification Module generates warning messages and security notifications whenever malicious activities are detected by the system. This module helps administrators and security analysts take immediate action against cyber threats before significant damage occurs.

### 10.6. Database Management Module

The Database Management Module is responsible for storing and managing all system-related information. This module stores datasets, user records, network traffic logs, machine learning prediction results, threat logs, alerts, and dashboard reports.

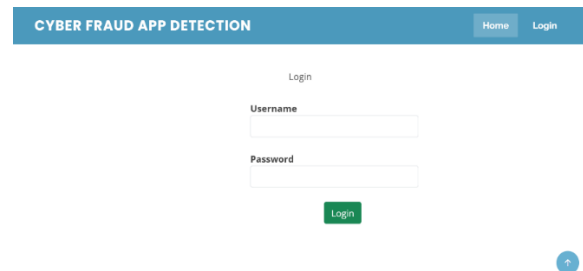




Fig 3: Login Page

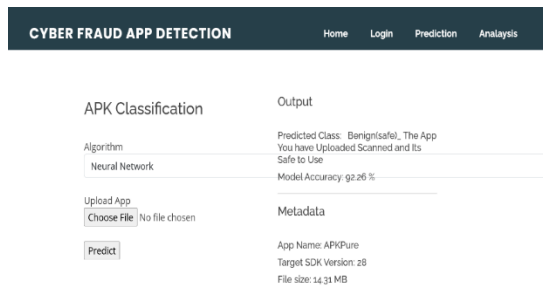


Fig 3: Detection Page

## 11. IMPLEMENTATION

The implementation of the “Cyber Threat Detection Using Machine Learning and Tableau” project involves the development of an intelligent cybersecurity system capable of detecting malicious activities in network traffic using machine learning algorithms and visualizing the results through interactive dashboards. The implementation process includes data collection, preprocessing, feature extraction, machine learning model training, threat detection, database management, and dashboard visualization. The project is implemented using Python programming language along with machine learning libraries and Tableau for graphical analysis.

The implementation process begins with collecting cybersecurity datasets and network traffic information from different sources. Publicly available datasets such as NSL-KDD, CICIDS2017, and UNSW-NB15 are used for training and testing the machine learning models. These datasets contain information about normal network behavior and different types of cyber-attacks such as malware, phishing, denial-of-service attacks, and unauthorized access attempts. The collected datasets are stored in CSV format and imported into the Python environment for further processing.

After data collection, the next implementation stage focuses on data preprocessing. The raw dataset often contains missing values, duplicate entries, irrelevant features, and inconsistent data formats that can reduce the performance of machine learning algorithms. Therefore, preprocessing techniques such as data cleaning, normalization, encoding categorical values, and feature scaling are applied. Duplicate records are removed, null values are handled, and the dataset is converted into a machine-readable format. This stage improves the quality and reliability of the dataset and helps achieve better prediction accuracy.

Once preprocessing is completed, feature extraction and feature selection techniques are implemented to identify important attributes required for cyber threat detection. Features such as source IP address, destination IP address, packet size, protocol type, traffic duration, and connection status are selected from the dataset. Redundant and unnecessary features are removed to reduce complexity and improve model efficiency. The selected features are then divided into training and testing datasets using train-test split techniques.

Overall, the implementation of this project demonstrates how machine learning and data analytics can be effectively applied in the field of cybersecurity to develop automated, scalable, and real-time threat detection solutions.

## 12. ALGORITHM USED

### 12.1 Neural Network Algorithm

A Neural Network is a machine learning algorithm inspired by the structure and functioning of the human brain. It consists of interconnected nodes called neurons that work together to process information and recognize patterns from large datasets. Neural Networks are highly effective in solving complex classification and prediction problems because they can learn hidden relationships between input data and output results.

In the cyber threat detection system, the Neural Network algorithm is used to analyze network traffic patterns and identify malicious activities. The algorithm receives network-related features such as source IP address, destination IP address, protocol type, packet size, traffic duration, and connection behavior as input data. These inputs are processed through multiple hidden layers where mathematical calculations and activation functions are applied to detect complex relationships within the data.

### 12.2 Support Vector Classifier (SVC)

Support Vector Classifier (SVC) is a supervised machine learning algorithm based on the concept of Support Vector Machines (SVM). It is mainly used for classification tasks and is highly effective for binary classification problems such as identifying whether network traffic is normal or malicious. The primary objective of the SVC algorithm is to find the optimal decision boundary, called a hyperplane, that separates different classes of data with maximum margin.

In the cyber threat detection project, the SVC algorithm analyzes cybersecurity data and classifies



network activities into safe or harmful categories. The algorithm works by mapping input data points into a multidimensional feature space and identifying the best possible boundary between normal and malicious activities. The data points closest to the hyperplane are called support vectors, and they play an important role in defining the classification boundary.

The Support Vector Classifier is highly effective in handling high-dimensional cybersecurity datasets where multiple features are involved. It can efficiently detect intrusion attempts, malware traffic, phishing activities, and unauthorized access patterns. SVC performs well even when the dataset contains complex relationships and overlapping classes.

### 13. RESULT AND DISCUSSION

The “Cyber Threat Detection Using Machine Learning ” project was successfully implemented and tested using various machine learning algorithms to detect malicious activities in network traffic data. The system was trained using cybersecurity datasets containing both normal and malicious network activities. Different machine learning algorithms such as Neural Network and Support Vector Classifier (SVC) were applied to classify cyber threats and evaluate the performance of the proposed system. The implementation demonstrated that machine learning techniques can effectively improve cyber threat detection accuracy and reduce manual monitoring efforts.

During the implementation process, the collected cybersecurity dataset was preprocessed to remove missing values, duplicate records, and irrelevant information. Feature extraction and feature selection techniques were applied to select important attributes such as protocol type, packet size, source IP address, destination IP address, connection duration, and traffic behavior. These optimized features improved the efficiency and accuracy of the machine learning models.

The Support Vector Classifier (SVC) algorithm also produced effective results in classifying network traffic into normal and malicious categories. SVC created an optimal decision boundary between normal and suspicious activities and performed well for binary classification tasks. The algorithm showed strong accuracy and reliable performance for intrusion detection and anomaly analysis. It also reduced classification errors and false positive rates compared to traditional rule-based security systems.

The project also integrated Tableau dashboards for cybersecurity visualization and reporting. The dashboard displayed attack statistics, traffic analysis, suspicious activities, model performance, and threat categories through graphs and charts. The visualization component improved system usability by helping administrators monitor cyber threats in real time and understand network behavior more effectively.

The overall discussion of the project indicates that machine learning algorithms provide significant advantages over traditional cybersecurity systems. Traditional systems mainly depend on predefined signatures and static rules, which are unable to detect unknown and evolving cyber threats efficiently. In contrast, the proposed machine learning-based system can learn attack patterns automatically and identify abnormal behavior with improved accuracy and speed. The integration of visualization tools further enhances cybersecurity analysis and decision-making capabilities.

### 14. FUTURE ENHANCEMENTS

One of the major future enhancements of the project is the implementation of Deep Learning techniques such as Convolutional Neural Networks (CNN), Recurrent Neural Networks (RNN), and Long Short-Term Memory (LSTM) models. Deep learning algorithms can analyze complex and high-dimensional cybersecurity data more effectively than traditional machine learning models. These techniques can improve the detection of advanced persistent threats, zero-day attacks, and hidden malicious activities by learning deep behavioral patterns from network traffic data.

Another important enhancement is the integration of real-time cyber threat monitoring and live network traffic analysis. In the current system, datasets are mainly used for offline training and testing purposes. Future versions of the project can be connected to real-time network environments where live traffic packets are continuously monitored and analyzed. This enhancement will help organizations detect cyber threats instantly and respond to attacks more quickly.

The project can also be enhanced by integrating cloud computing and distributed systems for large-scale cybersecurity analysis. Cloud-based implementation will improve scalability, storage capacity, and computational performance while handling massive volumes of network traffic data. Cloud integration can also support remote monitoring and centralized cybersecurity



management for organizations operating in multiple locations.

Another enhancement involves integrating advanced anomaly detection techniques and hybrid machine learning models. Combining supervised and unsupervised learning methods can improve the detection of unknown and evolving cyber threats. Hybrid models can also reduce false positive rates and improve the overall accuracy of the cybersecurity system.

The project can further be enhanced by supporting Internet of Things (IoT) security monitoring. As IoT devices become more common in industries, homes, and healthcare systems, cyber threats targeting IoT environments are increasing rapidly. Future versions of the project can analyze IoT network traffic and detect vulnerabilities or attacks specifically targeting connected smart devices.

Future development may also include mobile application integration for real-time monitoring and alert management. Administrators and security analysts could receive instant notifications and monitor cybersecurity dashboards directly through mobile devices, improving accessibility and response efficiency.

The visualization component using Tableau can also be improved by implementing more interactive dashboards, predictive analytics, real-time attack maps, and dynamic threat intelligence reports. Advanced visualizations will help organizations better understand attack patterns and security trends.

## 15. REFERENCES

- [1] Ian Goodfellow, Yoshua Bengio, and Aaron Courville, *Deep Learning*, MIT Press, 2016.
- [2] Stuart Russell and Peter Norvig, *Artificial Intelligence: A Modern Approach*, Pearson Education, 2010.
- [3] William Stallings, *Network Security Essentials: Applications and Standards*, Pearson Education, 2017.
- [4] Behrouz A. Forouzan, *Cryptography and Network Security*, McGraw Hill Education, 2015.
- [5] Nina Godbole and Sunit Belapure, *Cyber Security: Understanding Cyber Crimes, Computer Forensics and Legal Perspectives*, Wiley India, 2011.
- [6] Tom M. Mitchell, *Machine Learning*, McGraw Hill Education, 1997.
- [7] Ethem Alpaydin, *Introduction to Machine Learning*, MIT Press, 2020.
- [8] Trevor Hastie, Robert Tibshirani, and Jerome Friedman, *The Elements of Statistical Learning*, Springer, 2009.
- [9] Christopher M. Bishop, *Pattern Recognition and Machine Learning*, Springer, 2006.
- [10] Charu C. Aggarwal, *Neural Networks and Deep Learning*, Springer, 2018.