# PROJECT REPORT

## ON

# MALICIOUS DNS TRAFFIC DETECTION USING DEEP LEARNING



## CENTER FOR DEVELOPMENT OF ADVANCED COMPUTING

## ELECTRONIC CITY,BANGALORE

## UNDER THE SUPERVISION OF

### SUBHAM GOYAL

### SENIOR PROJECT ENGINEER

### C-DAC BANGALORE

### Submitted by:-

Abhishek Kumar Singh(240350125003)

Manchi.Tejasai(240350125037)

Ankit Kumar(240350125008)

Kothapalli Mounika(240350125035)

B.Dheeraj Chandan(240350120017)

## PG DIPLOMA IN ADVANCED COMPUTING/PG DIPLOMA IN IT
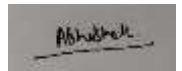
## INFRASTRUCTURE, SYSTEM & SECURITY

## C-DAC,BANGALORE
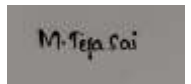
# Candidate Declaration

We hereby certify that the work being presented in the report entitled **MALICIOUS DNS TRAFFIC DETECTION USING DEEP LEARNING**, in the partial fulfilment of the requirements for the award of **PG-DIPLOMA** and submitted in the department of **DBDA** of the C-DAC Bangalore, is an authentic record of our work carried out during the period 1$^{st}$ July 2024– 14$^{th}$ August 2024 under the supervision of **SUBHAM GOYAL,SENIOR PROJECT ENGINEER**, C-DAC Bangalore.

The matter presented in the report has not been submitted by me for the award of any degree of this or any other Institute/University.
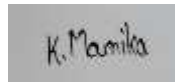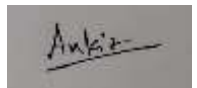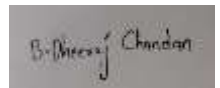
ABHISHEK KUMAR SINGH

MANCHI.TEJASAI

KOTHAPALLI MOUNIKA

ANKIT KUMAR

B.DHEERAJ CHANDAN

(Name and Signature of Candidate)

Counter Signed by

Name(s) of Supervisor(s)

....................................

....................................

# ACKNOWLEDGMENT

I take this opportunity to express my gratitude to all those people who have been directly and indirectly with me during the competition of this project

I pay thank to **Subham Goyal** who has given guidance and a light to me during this major project. His versatile knowledge about "title name" has eased me in the critical times during the span of this Final Project.

I acknowledge here out debt to those who contributed significantly to one or more steps. I take full responsibility for any remaining sins of omission and commission.

**Student Name**

ABHISHEK KUMAR SINGH

MANCHI.TEJASAI

KOTHAPALLI MOUNIKA

ANKIT KUMAR

B.DHEERAJ CHANDAN

# ABSTRACT

The project "DNS Traffic Detection Using Deep Learning" aims to enhance the detection of malicious activities within DNS traffic by leveraging advanced deep learning techniques. The Domain Name System (DNS) is a critical component of internet infrastructure, translating human-readable domain names into machine-readable IP addresses. However, its essential role also makes it a target for cyber threats such as DNS tunneling, amplification attacks, and domain generation algorithms (DGAs). Traditional detection methods, which rely on rule-based and signature-based systems, often suffer from high false positive rates, evasion techniques, and scalability issues.

This project explores the application of deep learning models, including Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and Autoencoders, to address these limitations. Deep learning offers automatic feature extraction, superior pattern recognition, and adaptability to evolving threats, making it well-suited for DNS traffic analysis. The project involves developing a robust deep learning model, comparing the performance of different architectures, and evaluating the model's effectiveness in real-time detection scenarios.

Key objectives include achieving high detection accuracy, ensuring model interpretability and scalability, and providing actionable insights for cybersecurity professionals. The project also identifies future directions for continuous improvement, such as hybrid models and integration with other security systems. Through this approach, the project seeks to advance the state-of-the-art in DNS traffic detection and contribute to a more secure internet environment.

# Table of content

# INTRODUCTION

## 1.1 Background on DNS

The Domain Name System (DNS) is a cornerstone of the internet's infrastructure, responsible for translating human-readable domain names (like [www.example.com](www.example.com)) into machine-readable IP addresses (like 192.0.2.1). This translation is essential for routing traffic across the internet, enabling users to access websites, send emails, and perform other online activities seamlessly.

DNS operates through a hierarchical structure of servers, starting from root servers down to authoritative servers for specific domains. When a user types a URL into their browser, a DNS query is sent to resolve the domain name into an IP address. This process involves multiple steps and interactions between various DNS servers.

## 1.2 Importance of DNS Traffic Monitoring

Given its critical role, DNS traffic is a valuable source of information for network administrators and cybersecurity professionals. Monitoring DNS traffic can provide insights into network performance, user behavior, and potential security threats. However, the same characteristics that make DNS essential also make it a target for malicious activities.

## 1.3 Threats Exploiting DNS

Cyber attackers exploit DNS in various ways, including:

- ***DNS Tunneling***: This technique involves encapsulating data within DNS queries and responses to bypass network security measures. Attackers can use DNS tunneling to exfiltrate sensitive data or establish command-and-control channels.

- ***DNS Amplification Attacks***: These are a type of Distributed Denial of Service (DDoS) attack where attackers send small DNS queries with spoofed IP addresses to open DNS resolvers, which then send large responses to the victim, overwhelming their network.

- ***Domain Generation Algorithms (DGAs):*** Malware often uses DGAs to generate a large number of domain names in an attempt to evade detection and ensure that at least some of the domains remain accessible for command-and-control purposes.

## 1.4 Limitations of Traditional Detection Methods

Traditional DNS traffic detection methods primarily rely on rule-based systems and signature detection. These methods have several limitations:

- ***High False Positive Rates***: Rule-based systems can generate a significant number of false positives, flagging legitimate traffic as malicious. This can overwhelm security teams and reduce the effectiveness of the detection system.

- ***Evasion Techniques***: Attackers continuously evolve their methods to bypass signature-based detection. For example, they may use encryption or obfuscation techniques to hide malicious activities within DNS traffic.

- ***Scalability Issues***: The sheer volume of DNS traffic in large networks can make it challenging to analyze in real-time using traditional methods. This can result in delayed detection and response to threats.

## 1.5 The Role of Deep Learning in DNS Traffic Detection

Deep learning, a subset of machine learning, involves training neural networks with multiple layers (hence "deep") to learn complex patterns from large datasets. Unlike traditional machine learning, which often requires manual feature extraction, deep learning models can automatically learn relevant features from raw data. This capability makes deep learning particularly well-suited for DNS traffic detection.

## 1.6 Advantages of Deep Learning for DNS Traffic Detection

- *Feature Learning*: Deep learning models can automatically learn and extract features from raw DNS traffic data, reducing the need for manual feature engineering. This allows the models to adapt to new types of attacks more effectively.

- *Pattern Recognition*: Neural networks excel at recognizing complex patterns in data. This makes them capable of identifying subtle indicators of malicious activity that traditional methods might miss.

- *Adaptability*: Deep learning models can be retrained with new data to adapt to evolving threats. This continuous learning process helps maintain high detection accuracy over time.

## 1.7 Deep Learning Architectures for DNS Traffic Detection

Several deep learning architectures can be employed for DNS traffic detection, each with its strengths:

- *Convolutional Neural Networks (CNNs)*: CNNs are effective for spatial data analysis and can be used to analyze sequences of DNS queries. They are particularly good at identifying local patterns in data.

- *Recurrent Neural Networks (RNNs)*: RNNs are designed for sequential data and can capture temporal dependencies in DNS traffic. This makes them suitable for detecting patterns that unfold over time, such as the behavior of DGAs.

- *Autoencoders*: Autoencoders are a type of unsupervised learning model used for anomaly detection. They learn a compressed representation of normal traffic and can identify deviations from this norm, flagging potential anomalies.

## 1.8 Project Goals

The primary goals of this project are to:

- **Develop a deep learning model capable of detecting malicious DNS traffic with high accuracy.**

- **Compare the performance of different deep learning architectures .**

- **Evaluate the model's effectiveness in real-time detection scenarios.**

- **Provide insights into the interpretability and scalability of the model.**

# Objective

## 1. Develop a Deep Learning Model to Detect Malicious DNS Traffic

Create a robust deep learning model that can accurately identify malicious DNS traffic.

**Details:**

- Model Selection: Choose appropriate deep learning architectures (e.g., CNNs, RNNs, Autoencoders) based on the nature of DNS traffic data.

- Feature Learning: Utilize the model's ability to automatically learn and extract relevant features from raw DNS traffic data, reducing the need for manual feature engineering.

- Training: Train the model using a comprehensive dataset that includes both benign and malicious DNS traffic. Ensure the dataset is balanced to prevent bias.

- Optimization: Fine-tune hyperparameters and optimize the model to achieve the best possible performance in terms of accuracy, precision, recall, and F1-score.

## 2. Compare the Performance of Different Deep Learning Architectures

Evaluate and compare the effectiveness of various deep learning architectures in detecting malicious DNS traffic.

**Details:**

- *Architectures:* Implement and test different deep learning architectures, such as:

  - Convolutional Neural Networks (CNNs): Effective for spatial data analysis and identifying local patterns in DNS queries.

  - Recurrent Neural Networks (RNNs): Suitable for sequential data analysis and capturing temporal dependencies in DNS traffic.

  - Autoencoders: Useful for anomaly detection by learning a compressed representation of normal traffic and identifying deviations.

- *Performance Metrics:* Use standard evaluation metrics (accuracy, precision, recall, F1-score) to compare the performance of each architecture.

- *Cross-Validation:* Employ cross-validation techniques to ensure the robustness and generalizability of the models.

## 3. Evaluate the Model's Effectiveness in Real-Time Detection Scenarios

Assess the model's ability to detect malicious DNS traffic in real-time environments.

**Details:**

- *Real-Time Processing:* Implement the model in a real-time processing pipeline to monitor DNS traffic as it occurs.

- *Latency:* Measure the model's processing latency to ensure it can operate within acceptable time frames for real-time detection.

- *Scalability:* Test the model's scalability by evaluating its performance on large volumes of DNS traffic, typical of real-world network environments.

- *False Positives/Negatives:* Analyze the rate of false positives and false negatives in real-time detection to ensure the model's reliability.

## 4. Provide Insights into the Interpretability and Scalability of the Model

Ensure that the deep learning model is interpretable and scalable, making it practical for deployment in real-world scenarios.

**Details:**

- *Model Interpretability:* Develop methods to interpret the model's decisions, providing actionable insights to security analysts. This may involve:

    - *Feature Importance:* Identifying which features the model considers most important for making decisions.

    - *Visualization:* Creating visualizations to help understand the model's behavior and decision-making process.

- *Scalability:* Ensure the model can handle large-scale DNS traffic without significant degradation in performance. This involves:

    - *Resource Efficiency:* Optimizing the model to run efficiently on available hardware resources.

    - *Distributed Processing:* Exploring distributed processing techniques to handle high volumes of DNS traffic.

- *Deployment:* Develop a deployment strategy that integrates the model into existing network security infrastructure, ensuring seamless operation and minimal disruption.

# Literature Review

## 1. Traditional DNS Traffic Analysis Methods

### 1.1 Rule-Based Systems

Rule-based systems rely on predefined rules and signatures to detect malicious DNS traffic. These rules are often based on known patterns of malicious behavior, such as specific domain names, query types, or response codes.

*Key Studies:*

- Bilge et al. (2011): Introduced EXPOSURE, a system that uses passive DNS analysis to detect malicious domains. The system relies on a set of heuristics and rules to identify suspicious patterns in DNS traffic.

- Yadav et al. (2010): Proposed a method to detect algorithmically generated domain names (DGAs) by analyzing the statistical properties of domain names and their query patterns.

*Limitations:*

- High False Positives: Rule-based systems often generate a significant number of false positives, flagging legitimate traffic as malicious.

- Evasion Techniques: Attackers can easily modify their tactics to bypass rule-based detection, rendering these systems less effective over time.

- Scalability: Managing and updating a large set of rules can be challenging, especially in dynamic and large-scale network environments.

### 1.2 Signature-Based Detection

Signature-based detection involves identifying known malicious patterns or signatures in DNS traffic. These signatures are typically derived from previous attacks and are stored in a database for comparison.

*Key Studies:*

- Antonakakis et al. (2011): Developed a system to detect malware domains at the upper DNS hierarchy by analyzing DNS query patterns and comparing them to known malicious signatures.

- Perdisci et al. (2009): Proposed a method to detect fast-flux service networks by identifying the rapid changes in IP addresses associated with malicious domains.

*Limitations:*

- Limited Scope: Signature-based detection is effective only against known threats. New or unknown attacks that do not match existing signatures can go undetected.

- Maintenance: Continuously updating the signature database to include new threats requires significant effort and resources.

## 2. Machine Learning Approaches

### 2.1 Supervised Learning

Supervised learning involves training a model on labeled data, where each data point is associated with a known outcome (e.g., benign or malicious). The model learns to classify new data based on the patterns observed in the training data.

*Key Studies:*

- Bilge et al. (2014): Proposed a machine learning-based approach to detect malicious domains by analyzing features such as domain name length, entropy, and query frequency.

- Antonakakis et al. (2012): Developed Pleiades, a system that uses supervised learning to detect DGAs by analyzing the lexical and temporal features of domain names.

*Limitations:*

- Feature Engineering: Supervised learning models often require extensive feature engineering to identify relevant characteristics of DNS traffic. This process can be time-consuming and may not capture all relevant features.

- Data Quality: The effectiveness of supervised learning models depends on the quality and representativeness of the labeled training data.

## 2.2 Unsupervised Learning

Unsupervised learning involves training a model on unlabeled data to identify patterns and anomalies. This approach is useful for detecting new or unknown threats that do not have labeled examples.

*Key Studies:*

- Perdisci et al. (2010): Proposed an unsupervised learning method to detect botnet command-and-control domains by clustering DNS query patterns and identifying outliers.

- Yadav et al. (2012): Developed an unsupervised anomaly detection system to identify malicious domains by analyzing the statistical properties of DNS traffic.

*Limitations:*

- Interpretability: Unsupervised learning models can be difficult to interpret, making it challenging to understand the reasons behind their decisions.

- False Positives: These models may generate false positives by flagging legitimate but unusual traffic as malicious.

# 3. Deep Learning Approaches

## 3.1 Convolutional Neural Networks (CNNs)

CNNs are a type of deep learning model that excels at identifying spatial patterns in data. They are particularly effective for analyzing sequences of DNS queries and detecting local patterns indicative of malicious activity.

*Key Studies:*

- Saxe and Berlin (2017): Demonstrated the effectiveness of CNNs in detecting malware by analyzing the byte sequences of executable files. This approach can be adapted to analyze DNS query sequences.

- Yu et al. (2018): Proposed a CNN-based method to detect malicious domains by analyzing the character-level patterns in domain names.

*Advantages:*

- Automatic Feature Extraction: CNNs can automatically learn relevant features from raw data, reducing the need for manual feature engineering.
- Pattern Recognition: CNNs are effective at recognizing complex patterns in data, making them suitable for detecting subtle indicators of malicious activity.

## 3.2 Recurrent Neural Networks (RNNs)

RNNs are designed for sequential data and can capture temporal dependencies in DNS traffic. This makes them suitable for detecting patterns that unfold over time, such as the behavior of DGAs.

*Key Studies:*

- Liu et al. (2019): Developed an RNN-based model to detect malicious DNS traffic by analyzing the temporal patterns of DNS queries.
- Huang et al. (2020): Proposed a Long Short-Term Memory (LSTM) network, a type of RNN, to detect DNS tunneling by analyzing the sequence of DNS queries and responses.

*Advantages:*

- Temporal Dependencies: RNNs can capture the temporal relationships between DNS queries, making them effective for detecting sequential anomalies.
- Adaptability: RNNs can adapt to new types of attacks by learning from the temporal patterns in the data.

## 3.3 Autoencoders

Autoencoders are a type of unsupervised learning model used for anomaly detection. They learn a compressed representation of normal traffic and can identify deviations from this norm, flagging potential anomalies.

*Key Studies:*

- Mirsky et al. (2018): Developed Kitsune, an autoencoder-based network intrusion detection system that can detect anomalies in network traffic, including DNS traffic.
- Kim et al. (2019): Proposed an autoencoder-based method to detect DNS tunneling by learning the normal patterns of DNS traffic and identifying deviations.

*Advantages:*

- Anomaly Detection: Autoencoders are effective at identifying deviations from normal traffic patterns, making them suitable for detecting unknown or novel threats.
- Unsupervised Learning: Autoencoders do not require labeled data, making them useful for detecting new types of attacks.

# System Design

## 3.1 System Design :

A design document is a comprehensive plan that outlines the architecture, components, and implementation details of a project, system, or feature. It serves as a guide for developers, stakeholders, and other team members throughout the development process.

This diagram represents a flowchart for a process, likely related to detecting anomalies or threats based on DNS logs. Here's a breakdown of each step in the process:

1. ***Data Collection***

   - **Action:** Collect DNS Logs

   - **Purpose:** Gather DNS logs from various sources. These logs contain information about DNS queries made within a network, useful for detecting suspicious activities or patterns.

2. ***Data Preprocessing***

   - **Action:** Clean & Preprocess

   - **Purpose:** Clean and preprocess the collected data by removing noise, handling missing values, and transforming the data into a suitable format for analysis.

3. ***Feature Extraction***

   - **Action:** Extract Features

   - **Purpose:** Identify and extract relevant features from the preprocessed data, such as query patterns, frequency of requests, and types of queries.

4. ***Model Training***

   - **Action:** Train Model

   - **Purpose:** Use the extracted features to train a machine learning model to identify patterns or anomalies in DNS traffic.
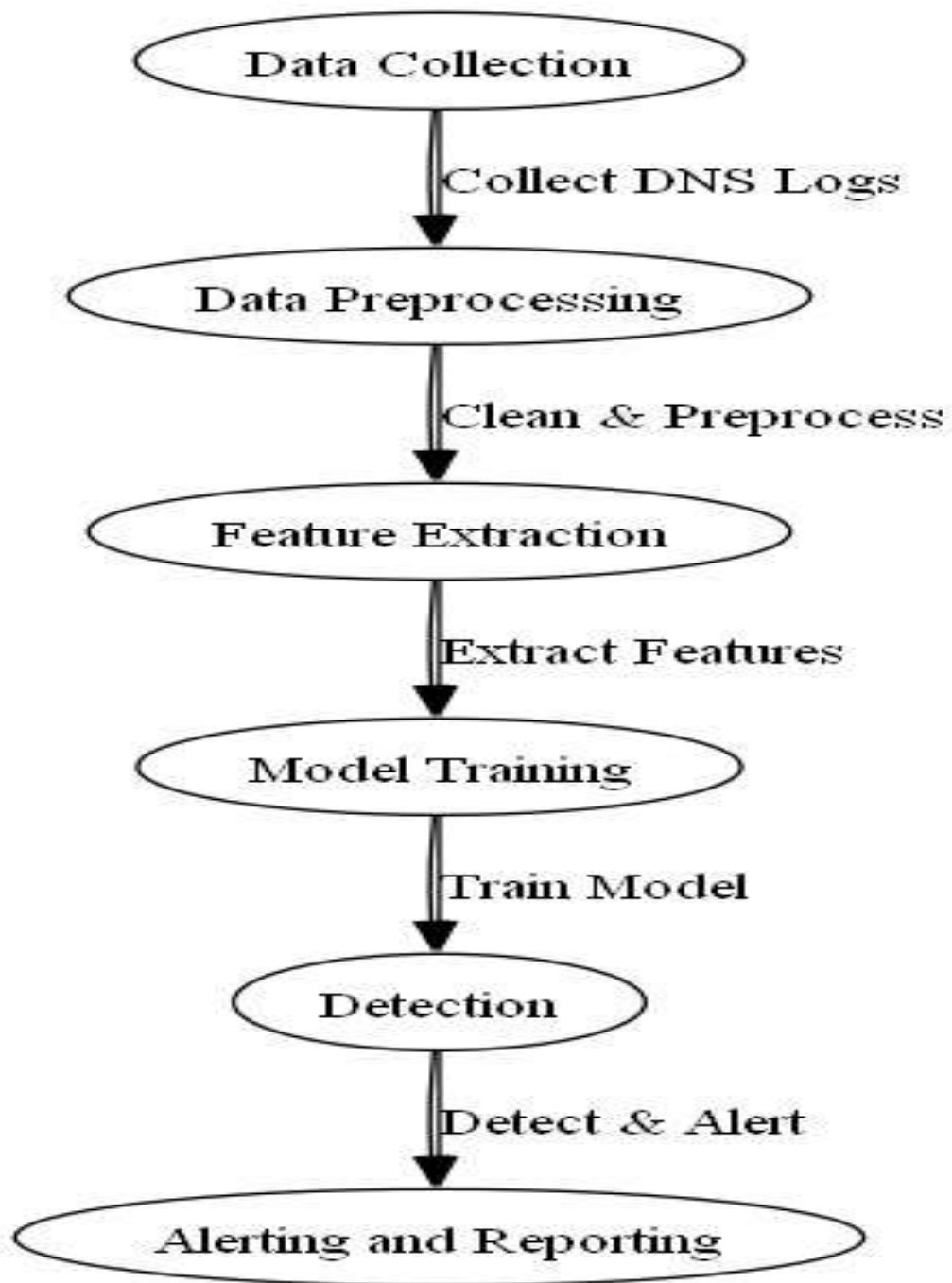
5. ***Detection***

   - **Action:** Detect & Alert

   - **Purpose:** Apply the trained model to new DNS logs to detect potential threats or anomalies and generate alerts.

6. ***Alerting and Reporting***

   - **Action:** Alert & Report

   - **Purpose:** Alert appropriate parties and generate reports based on detected anomalies, including details about the nature of the anomaly, its severity, and recommended actions.

This flowchart likely represents a high-level view of a system designed for monitoring and detecting suspicious activities in DNS logs, potentially for cybersecurity purposes. Each step in the process is critical for ensuring that the system can effectively identify and respond to threats in real-time.

Data Collection

Collect DNS Logs

Data Preprocessing

Clean & Preprocess

Feature Extraction

Extract Features

Model Training

Train Model

Detection

Detect & Alert

Alerting and Reporting

+

## 3.2 Enhanced Low-Level Diagram

This diagram is a flowchart that outlines the architecture of a data processing and analysis pipeline, particularly suited for time-series data and deep learning applications. Each module or component in the flowchart plays a specific role in handling, processing, and analyzing data from its raw form to producing real-time insights.

Here's a breakdown of each component:

1. ***Data Ingestion Module***

   - **Action:** Log Collection, Preprocessing

   - **Details:** Collect raw data, such as logs from various sources, and perform basic preprocessing.

2. ***Data Cleaning and Transformation Component***

   - **Action:** Data Cleaning, Data Transformation

   - **Details:** Clean the data to remove noise and inconsistencies, and transform it into a structured format.

3. ***Data Preprocessing***

   - **Action:** Normalization, Feature Extraction

   - **Details:** Normalize the data and extract important features for training machine learning models.

4. ***Data Storage Module***

   - **Action:** Time Series Database (InfluxDB), Data Management

   - **Details:** Store cleaned and preprocessed data in a time-series database optimized for handling data points indexed over time.

5. ***Data Retrieval***

   - **Action:** Querying TSDB (Time Series Database)

   - **Details:** Retrieve data from the time-series database for model training or real-time analysis.

6. ***Model Training Module***

   - **Action:** Deep Learning Models (TensorFlow, PyTorch), Training, Validation

   - **Details:** Train deep learning models using frameworks like TensorFlow or PyTorch, and validate their performance.

7. ***Model Evaluation***

   - **Action:** Performance Metrics, Comparison

   - **Details:** Evaluate models based on performance metrics to determine their effectiveness.

8. ***Model Deployment***

   - **Action:** Real-time Analysis

   - **Details:** Deploy the model for real-time analysis to provide predictions or insights.

9. **User Interface Module**

- **Action:** Web Dashboard, Visualization, Reporting

- **Details:** Provide visualization and reporting tools through a web dashboard for user interaction and insights.

*Summary:*

This flowchart outlines a comprehensive data processing pipeline from raw data collection to deploying a deep learning model for real-time analysis, and finally, presenting the results through a user-friendly interface. Each step in the process is crucial for ensuring the data is properly handled, transformed, and utilized to extract meaningful insights.

```
┌─────────────────────────────────────────┐
│         Data Ingestion Module            │
│     (Log Collection, Preprocessing)      │
└─────────────────────────────────────────┘
                    │ Data Ingestion
                    ▼
┌─────────────────────────────────────────┐
│ Data Cleaning and Transformation Component│
│  (Data Cleaning, Data Transformation)    │
└─────────────────────────────────────────┘
                    │ Data Cleaning and Transformation
                    ▼
┌─────────────────────────────────────────┐
│          Data Preprocessing              │
│   (Normalization, Feature Extraction)    │
└─────────────────────────────────────────┘
                    │ Data Preprocessing
                    ▼
┌─────────────────────────────────────────┐
│         Data Storage Module              │
│ (Time Series Database (InfluxDB), Data Management)│
└─────────────────────────────────────────┘
                    │ Data Storage
                    ▼
┌─────────────────────────────────────────┐
│          Data Retrieval                  │
│         (Querying TSDB)                   │
└─────────────────────────────────────────┘
                    │ Data Retrieval
                    ▼
┌─────────────────────────────────────────┐
│          Model Training Module           │
│(Deep Learning Models (TensorFlow, PyTorch), Training, Validation)│
└─────────────────────────────────────────┘
                    │ Model Training
                    ▼
┌─────────────────────────────────────────┐
│          Model Evaluation                │
│  (Performance Metrics, Comparison)       │
└─────────────────────────────────────────┘
                    │ Model Evaluation
                    ▼
┌─────────────────────────────────────────┐
│          Model Deployment                │
│        (Real-time Analysis)              │
└─────────────────────────────────────────┘
                    │ Model Deployment
                    ▼
┌─────────────────────────────────────────┐
│         User Interface Module            │
│ (Web Dashboard, Visualization, Reporting)│
└─────────────────────────────────────────┘
```

## 3.3 Response code :

This diagram illustrates the structure of DNS (Domain Name System) logs by breaking down the various components that these logs typically contain. DNS logs are critical in network management and cybersecurity, as they provide detailed records of DNS queries and responses, which can be used for monitoring, troubleshooting, and detecting malicious activities.

Here's a breakdown of each element in the diagram:

### 1. DNS Logs :

**Description :** This is the main container or source of information in the diagram. DNS logs capture the details of DNS queries and responses made within a network.

### 2. Timestamp :

**Contains :** The specific date and time when the DNS query was made or when the DNS response was received. This information is crucial for tracking the timing of requests, which can be important for identifying patterns or incidents.

### 3. Source IP :

**Contains :** The IP address of the device that made the DNS query. This helps in identifying which device on the network initiated the DNS request.

### 4. Destination IP :

**Contains :** The IP address of the DNS server that received the query. This shows where the DNS query was sent to resolve the domain name.

### 5. Query Name :

**Contains :** The domain name that was queried. This is the specific domain name that the device was attempting to resolve into an IP address.

### 6. Query Type :

**Contains :** The type of DNS query made. Common query types include `A` (address record), `AAAA` (IPv6 address record), `MX` (mail exchange record), `CNAME` (canonical name record), etc. This indicates the type of information the client was requesting.

### 7. Response Code :

***Contains :*** The response code returned by the DNS server. This code indicates the result of the DNS query. Common response codes include `NOERROR` (query was successful), `NXDOMAIN` (non-existent domain), `SERVFAIL` (server failure), etc. The response code is important for diagnosing issues with DNS queries.

***Summary :***

This diagram provides a detailed view of the typical contents of DNS logs, showing that each log entry includes a timestamp, source and destination IPs, the queried domain name, the type of query, and the response code. Understanding these components is essential for network administrators and cybersecurity professionals to monitor DNS activity and troubleshoot any potential issues.

# Methodology

The methodology section outlines the systematic approach taken to develop and evaluate a deep learning model for detecting malicious DNS traffic. This involves several key steps, including data collection, data preprocessing, model selection, training, and evaluation. Each step is crucial for ensuring the effectiveness and reliability of the final model.

## 4.1 Data Collection

Data collection is a critical step in developing a deep learning model for DNS traffic detection. The quality and comprehensiveness of the collected data directly impact the model's performance and its ability to generalize to real-world scenarios. Here's a detailed breakdown of the data collection process:

### 4.1.1 Sources of Data

To build a robust dataset, it is essential to gather DNS traffic data from diverse and reliable sources. This ensures that the dataset captures a wide range of normal and malicious behaviors.

### 4.1.1.1 Public Datasets

Public datasets are often curated by research institutions, cybersecurity organizations, and academic projects. These datasets are valuable because they are typically well-documented and include a variety of DNS traffic patterns.

- **Examples**:
  - CAIDA (Center for Applied Internet Data Analysis) : Provides anonymized DNS traffic datasets collected from various network vantage points.
  - ISC (Internet Systems Consortium) : Offers DNS query data from root servers and other DNS infrastructure.
  - DNSDB (DNS Database) : A passive DNS replication system that collects and stores DNS query and response data.

- **Advantages**:
  - Diversity: Public datasets often include traffic from different networks and regions, providing a broad view of DNS behavior.
  - Documentation: These datasets are usually well-documented, making it easier to understand the data structure and context.

### 4.1.1.2 Simulated Environments

Simulated environments allow researchers to generate synthetic DNS traffic under controlled conditions.

- **Examples**:
  - Custom DNS Servers: Set up DNS servers to generate and log DNS queries and responses.
  - Network Simulation Tools: Use tools like Mininet or GNS3 to simulate network environments and generate DNS traffic.

- **Advantages**:

- Controlled Conditions: Allows for precise control over the types of traffic generated, including specific attack patterns.
- Reproducibility: Simulated environments can be easily reproduced and modified to test different scenarios.

### 4.1.1.3 Network Logs

Collecting DNS traffic logs from real-world networks provides a realistic view of DNS behavior in operational environments.

- **Examples**:
  - Enterprise Networks: Collect DNS logs from corporate networks, including internal and external DNS queries.
  - ISP Networks: Obtain DNS traffic data from internet service providers, capturing a wide range of user behaviors.

- **Advantages**:
  - Realism: Real-world network logs provide a realistic view of DNS traffic, including legitimate and malicious activities.
  - Volume: Large volumes of data can be collected, providing a rich dataset for training and evaluation.

## 4.1.2 Types of Data

To train a deep learning model effectively, it is essential to include both benign and malicious DNS traffic in the dataset.

### 4.1.2.1 Benign Traffic

Benign traffic represents normal DNS queries and responses from legitimate users and applications.

- **Examples**:
  - Web Browsing: DNS queries generated by users accessing websites.
  - Email Communication: DNS queries related to email servers (e.g., MX records).
  - Software Updates: DNS queries from applications checking for updates.

- **Importance**:
  - Baseline: Provides a baseline of normal DNS behavior, which is essential for identifying deviations indicative of malicious activity.
  - Training: Helps the model learn the characteristics of legitimate traffic, reducing false positives.

### 4.1.2.2 Malicious Traffic

Malicious traffic includes DNS queries and responses associated with known malicious activities.

- **Examples**:
  - DNS Tunneling: DNS queries used to tunnel other protocols and exfiltrate data.
  - DGA-Generated Domains: Domains generated by malware using domain generation algorithms to evade detection.

- DNS Amplification Attacks: DNS queries designed to exploit open resolvers and amplify attack traffic.

- **Importance**:

  - Detection: Provides examples of malicious behavior for the model to learn and detect.

  - Evaluation: Allows for the evaluation of the model's ability to identify different types of attacks.

## 4.1.3 Labeling

Labeling the data is a crucial step in supervised learning, as it provides the ground truth for training and evaluating the model.

### 4.1.3.1 Manual Labeling

Security experts manually review and label the data based on known malicious domains and patterns.

- **Process**:

  - Domain Reputation: Check the reputation of domains using threat intelligence feeds and blacklists.

  - Pattern Analysis: Analyze query patterns, such as frequency and timing, to identify suspicious behavior.

  - Expert Review: Security analysts review the data and assign labels (benign or malicious) based on their expertise.

- **Advantages**:

  - Accuracy: Manual labeling by experts ensures high accuracy and reliability of the labels.

  - Context: Experts can consider the context and nuances of the data, improving the quality of the labels.

### 4.1.3.2 Automated Labeling

Automated tools and threat intelligence feeds can be used to label the data based on known malicious domains and patterns.

- **Process**:

  - Threat Intelligence Feeds: Use feeds from cybersecurity organizations to identify known malicious domains.

  - Automated Tools: Employ tools that analyze DNS traffic and automatically label suspicious queries based on predefined rules and heuristics.

- **Advantages**:

  - Scalability: Automated labeling can handle large datasets efficiently, reducing the time and effort required for manual labeling.

  - Consistency: Automated tools provide consistent labeling, minimizing human error and bias.

## 4.2 Data Preprocessing

Data preprocessing is a crucial step in preparing raw DNS traffic data for training deep learning models. It involves cleaning, normalizing, and transforming the data to ensure it is in a suitable format for model training. Proper preprocessing helps improve the model's performance and generalizability. Here's a detailed breakdown of the data preprocessing process:

### 4.2.1 Data Cleaning

Remove noise, inconsistencies, and irrelevant information from the raw data to ensure its quality and reliability.

- **Details**:
    - Remove Incomplete Records
    - Filter Noise
    - Deduplication

### 4.2.2 Data Normalization

Scale numerical features to a standard range and encode categorical features to ensure uniformity and improve model performance.

- **Details**:
    - Scale Numerical Features: Normalize numerical features (e.g., query length, response time) to a standard range (e.g., 0 to 1) to ensure uniformity.
    - Encode Categorical Features: Convert categorical features (e.g., query type, response code) into numerical representations using techniques like one-hot encoding.

### 4.2.3 Feature Extraction

Extract relevant features from the raw DNS traffic data that can help the model distinguish between benign and malicious traffic.

- **Details**:
    - Query Length: Length of the DNS query string.
    - Query Type: Type of DNS query (e.g., A, AAAA, MX, TXT).
    - Response Time: Time taken to receive a response for the DNS query.
    - Domain Entropy: Measure of randomness in the domain name, useful for detecting DGAs.
    - Frequency of Queries: Number of queries for a particular domain within a specific time frame.
    - Time of Day: The time at which the query was made, which can help identify unusual patterns.
    - Geolocation: The geographic location of the DNS server or client, which can help identify unusual access patterns.
    -

### 4.2.4 Data Transformation

Transform the data into a format suitable for input into deep learning models.

- **Details**:

- Sequence Generation: For models like RNNs that require sequential data, generate sequences of DNS queries.

- Windowing: Divide the data into fixed-size windows to capture temporal patterns.

- Padding and Truncation: Ensure that all sequences have the same length by padding shorter sequences and truncating longer ones.

### 4.2.5 Data Splitting

Split the preprocessed data into training, validation, and test sets to ensure robust model evaluation.

- **Details**:

  - Training Set: Used to train the model.

  - Validation Set: Used to tune hyperparameters and prevent overfitting.

  - Test Set: Used to evaluate the final model's performance.

## 4.3 Model Selection

Model selection is a critical step in developing a deep learning-based DNS traffic detection system. The choice of model architecture significantly impacts the system's ability to accurately detect malicious activities while maintaining efficiency and scalability. This section provides a detailed explanation of the different deep learning architectures considered for DNS traffic detection, their strengths, and their suitability for the task.

### 4.3.1 Convolutional Neural Networks (CNNs)

CNNs are a type of deep learning model that excels at identifying spatial patterns in data.

- **Architecture**:

  - Convolutional Layers: Apply convolutional filters to the input data to extract local features.

  - Pooling Layers: Reduce the dimensionality of the data by down-sampling, which helps in capturing the most important features while reducing computational complexity.

  - Fully Connected Layers: Combine the features extracted by the convolutional layers to make the final classification.

- **Application to DNS Traffic**:

  - Input Representation: DNS queries can be represented as sequences of characters or tokens, which are then fed into the convolutional layers.

  - Pattern Recognition: CNNs can identify patterns such as unusual domain name structures, query frequencies, and response times that may indicate malicious activity.

- **Advantages**:

  - Automatic Feature Extraction: CNNs can automatically learn relevant features from raw data, reducing the need for manual feature engineering.

  - Pattern Recognition: Effective at recognizing complex patterns in data, making them suitable for detecting subtle indicators of malicious activity.

### 4.3.2 Recurrent Neural Networks (RNNs)

RNNs are designed for sequential data and can capture temporal dependencies in DNS traffic.

- **Architecture**:

- Recurrent Layers: Use recurrent units (e.g., LSTM or GRU) to process sequences of data, maintaining a hidden state that captures information from previous time steps.

  - Fully Connected Layers: Combine the features extracted by the recurrent layers to make the final classification.

- **Application to DNS Traffic**:

  - Input Representation: DNS queries are represented as sequences of characters or tokens, which are then fed into the recurrent layers.

  - Temporal Dependencies: RNNs can capture the temporal relationships between DNS queries, making them effective for detecting sequential anomalies.

- **Advantages**:

  - Temporal Dependencies: RNNs can capture the temporal relationships between DNS queries, making them effective for detecting sequential anomalies.

  - Adaptability: RNNs can adapt to new types of attacks by learning from the temporal patterns in the data.

### 4.3.3 Autoencoders

Autoencoders are a type of unsupervised learning model used for anomaly detection.

- **Architecture**:

  - Encoder: Compresses the input data into a lower-dimensional representation.

  - Decoder: Reconstructs the input data from the compressed representation.

  - Anomaly Detection: Anomalies are detected based on the reconstruction error, with higher errors indicating potential anomalies.

- **Application to DNS Traffic**:

  - Input Representation: DNS queries are represented as sequences of characters or tokens, which are then fed into the encoder.

  - Anomaly Detection: The model learns the normal patterns of DNS traffic and identifies deviations based on the reconstruction error.

- **Advantages**:

  - Anomaly Detection: Autoencoders are effective at identifying deviations from normal traffic patterns, making them suitable for detecting unknown or novel threats.

  - Unsupervised Learning: Autoencoders do not require labeled data, making them useful for detecting new types of attacks.

### 4.3.4 Model Selection Criteria

Choose the most suitable model architecture based on the specific requirements and characteristics of the DNS traffic detection task.

- **Details**:

  - Accuracy: Evaluate the model's ability to correctly classify benign and malicious DNS traffic.

- Precision and Recall: Assess the model's precision (proportion of true positives out of the total instances classified as positive) and recall (proportion of true positives out of the total actual positive instances).

- F1-Score: Use the F1-score, the harmonic mean of precision and recall, to provide a balanced measure of the model's performance.

- Scalability: Consider the model's ability to handle large volumes of DNS traffic in real-time.

- Computational Efficiency: Evaluate the model's computational requirements and ensure it can be deployed efficiently in a real-world environment.

- Interpretability: Ensure the model's decisions are interpretable and explainable to security analysts.

## 4.4 Training and Evaluation

Training and evaluation are crucial steps in developing a deep learning model for DNS traffic detection. These steps ensure that the model learns to accurately identify malicious DNS traffic and generalizes well to new, unseen data. Here's a detailed breakdown of the training and evaluation process:

### 4.4.1 Training

Train the selected deep learning model using the preprocessed data to learn patterns that distinguish between benign and malicious DNS traffic.

### 4.4.1.1 Dataset Split

Split the dataset into training, validation, and test sets to ensure robust model evaluation and prevent overfitting.

**Details:**

- Training Set: Used to train the model. Typically, 70% of the data is allocated to the training set.

- Validation Set: Used to tune hyperparameters and monitor the model's performance during training. Typically, 15% of the data is allocated to the validation set.

- Test Set: Used to evaluate the final model's performance. Typically, 15% of the data is allocated to the test set.

### 4.4.1.2 Balanced Dataset

Ensure a balanced representation of benign and malicious traffic in the training set to prevent bias.

**Details:**

- Class Imbalance: If the dataset is imbalanced (e.g., more benign traffic than malicious traffic), use techniques like oversampling, undersampling, or synthetic data generation (e.g., SMOTE) to balance the classes.

### 4.4.1.3 Hyperparameter Tuning

Optimize hyperparameters to improve the model's performance.

**Details:**

- Grid Search: Perform an exhaustive search over a specified parameter grid.

- Random Search: Perform a random search over a specified parameter grid.

- Bayesian Optimization: Use probabilistic models to find the optimal hyperparameters.

### *4.4.1.4 Regularization*

Apply regularization techniques to prevent overfitting.

**Details:**

- Dropout: Randomly drop units (along with their connections) during training to prevent overfitting.

- L2 Regularization: Add a penalty equal to the sum of the squared values of the weights to the loss function.

## 4.4.2 Evaluation

Evaluate the trained model's performance using standard metrics to ensure it generalizes well to new, unseen data.

### 4.4.2.1 Evaluation Metrics

Use standard evaluation metrics to assess the model's performance.

**Details:**

- Accuracy: Proportion of correctly classified instances out of the total instances.

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

- Precision: Proportion of true positive instances out of the total instances classified as positive.

  - Formula: $Precision = \frac{TP}{TP + FP}$

- Recall: Proportion of true positive instances out of the total actual positive instances.

  - Formula: $Recall = \frac{TP}{TP + FN}$

- F1-Score: Harmonic mean of precision and recall, providing a balanced measure of the model's performance.

  - Formula: $F1\text{-}Score = 2 \times \frac{Precision \times Recall}{Precision + Recall}$

- Confusion Matrix: A table that summarizes the performance of the model by showing the true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN).

### *4.4.2.2 Cross-Validation*

Use cross-validation techniques to ensure the model's robustness and generalizability.

**Details:**

- K-Fold Cross-Validation: Divide the dataset into K subsets and train the model K times, each time using a different subset as the validation set and the remaining subsets as the training set.

    - Example: 5-fold cross-validation involves splitting the data into 5 subsets and training the model 5 times.

### 4.4.2.3 Model Performance Monitoring

Continuously monitor the model's performance during training to detect overfitting and adjust training parameters as needed.

**Details:**

- Early Stopping: Stop training when the model's performance on the validation set stops improving.

- Learning Rate Scheduling: Adjust the learning rate during training to improve convergence.

# IMPLEMENTATION

## 1. Loading of csv:

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
df = pd.read_csv('data.csv', names=["Timestamp", "Duration", "Type", "Level",
                    "Client", "Client ID", "Query ID",
                    "Query Name", "View", "Recursion",
                    "Query Type", "Query", "Class",
                    "Record Type", "Flags", "IP Address"],low_memory=False)
df.head()
```

| | Timestamp | Duration | Type | Level | Client | Client ID | Query ID | Query Name | View | Recursion | Query Type | Query | Class | Record Type | Flags | IP |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 16-Apr-24 | 50:53.8 | queries: | info: | client | @0x7fccd4084398 | -4.013320e+18 | (fmcdg2.db-schenker.biz): | view | recursive: | query: | fmcdg2.db-schenker.biz | IN | A | + | (103.5 |
| 1 | 16-Apr-24 | 50:53.8 | queries: | info: | client | @0x7fcd2c4568a8 | -8.241670e+18 | (ip-10-0-44-250.us-west-2.compute.regio-bayeri... | view | recursive: | query: | ip-10-0-44-250.us-west-2.compute.regio-bayeris... | IN | A | + | (103.5 |
| 2 | 16-Apr-24 | 50:53.9 | queries: | info: | client | @0x7fcce41535b8 | -8.235850e+18 | (amazon.com): | view | recursive: | query: | amazon.com | IN | A | + | (103.5 |
| 3 | 16-Apr-24 | 50:53.9 | queries: | info: | client | @0x7fcce415d148 | -8.235850e+18 | (wikipedia.org): | view | recursive: | query: | wikipedia.org | IN | A | + | (103.5 |

## 2.Checking null value and row of data

```python
# Get the shape of the DataFrame
data_shape = df.shape

# Display the shape
print(f"Number of rows: {data_shape[0]}")
print(f"Number of columns: {data_shape[1]}")
```

```
Number of rows: 3692
Number of columns: 16
```

```python
null_values = df.isnull().sum()
print(null_values)
```

```
Timestamp      0
Duration       0
Type           0
Level          0
Client         0
Client ID      0
Query ID       0
Query Name     0
View           0
Recursion      0
Query Type     0
Query          0
Class          0
```

## 3.Displaying rows with null values and removing null values

```
# Display rows with any null values
rows_with_nulls = df[df.isnull().any(axis=1)]

# Display the rows with null values
print(rows_with_nulls)
```

```
Empty DataFrame
Columns: [Timestamp, Duration, Type, Level, Client, Client ID, Query ID, Query Name, View, Recursion, Query Type, Query, Class, Record Type, Flags, IP Address]
Index: []
```

```
[ ] # Drop rows with any null values
df.dropna(inplace=True)

# Verify that the rows with null values have been dropped
print(df.isnull().sum())
print(df.shape)
```

```
Timestamp      0
Duration       0
Type           0
Level          0
Client         0
Client ID      0
Query ID       0
Query Name     0
```

## 4.Choosing the column

```
(3692, 16)
```

```
# list of columns to select
columns_to_select = ["Timestamp", "Duration", "Client", "Client ID", "Query ID", "Query", "Class", "Record Type", "Flags", "IP Address"]

# Select the specified columns
selected_df = df[columns_to_select]

# Display the first few rows of the selected DataFrame
print(selected_df.head())

# Create a copy of the DataFrame
selected_columns = selected_df.copy()
selected_columns.shape
```

```
  Timestamp Duration Client        Client ID       Query ID \
0 16-Apr-24  50:53.8  client  @0x7fccd4084398  -4.013320e+18
1 16-Apr-24  50:53.8  client  @0x7fcd2c456Ba8  -8.241670e+18
2 16-Apr-24  50:53.9  client  @0x7fcce41535b8  -8.235850e+18
3 16-Apr-24  50:53.9  client  @0x7fcce415d148  -8.235850e+18
4 16-Apr-24  50:53.9  client  @0x7fcdd88263e8  -8.797800e+18

                                      Query Class Record Type Flags \
0                        fmcdg2.db-schenker.biz   IN           A     +
1  ip-10-0-44-250.us-west-2.compute.regin-bayoris...   IN           A     +
2                                    amazon.com   IN           A     +
3                                 wikipedia.org   IN           A     +
```

## 5.Converting time into proper format



## 6.Dropping the duration column

## 7.Making regex to pass the query



## 8.Checking percentage of query matched

## 9.Function to show unmatched queries



## 10. Adding regex function for unmatched queries

## 11.Again passing whole query column through regex

```
import re
import pandas as pd

query_column = 'Query'

# Get the unique unmatched queries
unmatched_queries = df[query_column].unique()

# Initialize the regex patterns dictionary
regex_patterns = {}

# Function to automatically update regex patterns
def update_regex_patterns(unique_queries, regex_patterns):
    for query in unique_queries:
        # Create a regex pattern for each unique unmatched query
        regex_patterns[f"unique_query_{query}"] = re.escape(query)
    return regex_patterns

# Update the regex patterns with unique unmatched queries
regex_patterns = update_regex_patterns(unmatched_queries, regex_patterns)

# Print the updated regex patterns
for name, pattern in regex_patterns.items():
    print(f"{name}: {pattern}")

# Create a new column in the DataFrame to store the matched patterns
df['matched_pattern'] = None

# Iterate over the rows in the DataFrame
for index, row in df.iterrows():
    query = row[query_column]
    matched_pattern = None
    # Iterate over the updated regex patterns to find a match
    for name, pattern in regex_patterns.items():
        if re.search(pattern, query):
            matched_pattern = name
```

## 12.Extracting IP address and saving it in txt

```
# Calculate the percentage of matched queries
num_matched_queries = df['matched_pattern'].notnull().sum()
total_queries = len(df)
percentage_matched_queries = (num_matched_queries / total_queries) * 100

# Print the percentage of matched queries
print(f"Percentage of matched queries: {percentage_matched_queries}%")
```

```
Percentage of matched queries: 100.0%
```

```
#Extracting the IP Address and saving it in txt file
# Function to extract IP addresses from the 'IP Address' column
def extract_ip_addresses(ip_string):
    # the regex to find all valid IPv4 and IPv6 addresses
    ip_pattern = r'(\d{1,3}\.\d{1,3}\.\d{1,3}\.\d{1,3})|([0-9a-fA-F:]+)'
    return re.findall(ip_pattern, ip_string)

# Extract IP addresses from the 'IP Address' column
ip_addresses = selected_columns['IP Address'].apply(extract_ip_addresses)

# Flatten the list of lists and get unique IP addresses
unique_ips = set(ip for sublist in ip_addresses for ip in sublist)

# Save the unique IP addresses to a text file
with open('IP_list.txt', 'w') as file:
    for ip in unique_ips:
        file.write(f"{ip}\n")  # Write each IP on a new line

print("Unique IP addresses saved to IP_list.txt")
```

```
Unique IP addresses saved to IP_list.txt
```

## 13.Saving domain in txt file



## 14.Percentage of domain extracted unique

**15.passing the domain through function to check which domain is benign or not**



**16. Loading the domain txt file**



**17.Merging the older file and the domain file**

## 18. Mapping the data

```python
# Define mappings for Category and Risk
category_mapping = {
    'Malicious': 1,
    'Suspicious': 0.5,  # Optional if you want to include this category
    'Benign': 0
}

risk_mapping = {
    'High': 2,
    'Medium': 1,
    'Low': 0
}

# Convert Category and Risk to numerical values
merged_df['Category'] = merged_df['Category'].replace(category_mapping)
merged_df['Risk'] = merged_df['Risk'].replace(risk_mapping)

# Display the updated DataFrame
print(merged_df[['Domain', 'Category', 'Risk']])
```

```
                                  Domain  Category  Risk
0                     fwcdg2.db-schenker.biz        0     0
1        compute.regio-bayerisch-schwaben.de        0     0
2        compute.regio-bayerisch-schwaben.de        0     0
3        compute.regio-bayerisch-schwaben.de        0     0
4                              2002205.1200        0     0
```

## 19 . Cleaning the ip address and checking how many time a query is called from one ip address

## 20.Label encoding the data



## 21. Using sequential model

## 22.Score of sequential



## 23.cnn model

## 24. cnn score



## 25.lstm model

## 26.lstm score

# MODEL INTERPRETABILITY AND SCALABILITY

## 6.1 Model Interpretability

Ensure that the deep learning model's decisions are interpretable and explainable to security analysts, providing actionable insights for threat detection and response.

### 6.1.1 Feature Importance

Identify which features the model considers most important for making decisions. This helps analysts understand the factors influencing the model's predictions.

*Techniques:*

- **Permutation Feature Importance:** Measure the change in the model's performance when a feature's values are randomly shuffled. A significant drop in performance indicates that the feature is important.

- **SHAP (SHapley Additive exPlanations):** A game-theoretic approach that assigns each feature an importance value for a particular prediction.

*Example:*

import shap

*# Assuming 'model' is a trained model and 'X_train' is the training data*

explainer = shap.Explainer(model, X_train)

shap_values = explainer(X_train)

*# Plot feature importance*

shap.summary_plot(shap_values, X_train)

### 6.1.2 Visualization

Create visualizations to help understand the model's behavior and decision-making process. Visualizations can make complex model outputs more accessible and interpretable.

*Techniques:*

- **Heatmaps:** Visualize the importance of different features or input regions.

- **Attention Maps:** For models with attention mechanisms, visualize which parts of the input the model focuses on when making a decision.

*Example:*

import seaborn as sns

import matplotlib.pyplot as plt

*# Assuming 'shap_values' is obtained from SHAP*

shap.summary_plot(shap_values, X_train, plot_type="bar")

*# Heatmap example*

sns.heatmap(shap_values.values, cmap='viridis')

plt.show()

### 6.1.3 Explainable AI (XAI)

Use techniques from the field of Explainable AI to provide explanations for the model's predictions. This helps build trust in the model and ensures that its decisions are transparent.

*Techniques:*

- **LIME (Local Interpretable Model-agnostic Explanations):** Explain individual predictions by approximating the model locally with an interpretable model.

- **Integrated Gradients:** Attribute the prediction of a deep network to its input features by integrating gradients along the path from a baseline input to the actual input.

*Example:*

import lime

from lime.lime_tabular import LimeTabularExplainer

*# Assuming 'model' is a trained model and 'X_train' is the training data*

explainer = LimeTabularExplainer(X_train, feature_names=feature_names, class_names=['benign', 'malicious'], mode='classification')

exp = explainer.explain_instance(X_test[0], model.predict_proba, num_features=10)

exp.show_in_notebook(show_table=True)

## 6.2 Model Scalability

Ensure that the deep learning model can handle large volumes of DNS traffic efficiently, making it practical for deployment in real-world network environments.

### 6.2.1 Resource Efficiency

Optimize the model to run efficiently on available hardware resources, ensuring that it can process DNS traffic in real-time without significant delays.

*Techniques:*

- **Model Compression:** Reduce the size of the model using techniques like pruning, quantization, and knowledge distillation.

- **Efficient Architectures:** Use lightweight model architectures designed for efficiency, such as MobileNet or EfficientNet.

*Example:*

from tensorflow_model_optimization.sparsity import keras as sparsity

*# Pruning example*

pruning_params = {

  'pruning_schedule': sparsity.PolynomialDecay(initial_sparsity=0.50, final_sparsity=0.90, begin_step=2000, end_step=10000)

}

model = sparsity.prune_low_magnitude(model, **pruning_params)

model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

### 6.2.2 Distributed Processing

Implement distributed processing frameworks to handle high volumes of DNS traffic, ensuring that the model can scale horizontally.

*Techniques:*

- **Apache Kafka:** Use Kafka for real-time data ingestion and streaming.

- **Apache Spark:** Use Spark for distributed data processing and model inference.

*Example:*

```
from pyspark.sql import SparkSession

from pyspark.ml import PipelineModel
```

*# Initialize Spark session*

```
spark = SparkSession.builder.appName("DNSDetection").getOrCreate()
```

*# Load pre-trained model*

```
model = PipelineModel.load("path/to/spark/model")
```

*# Read streaming data from Kafka*

```
df = spark.readStream.format("kafka").option("kafka.bootstrap.servers",
"localhost:9092").option("subscribe", "dns_topic").load()
```

*# Preprocess and make predictions*

```
predictions = model.transform(df)
```

### 6.2.3 Real-Time Processing

Implement the trained model in a real-time processing pipeline to monitor DNS traffic and detect anomalies as they occur.

*Techniques:*

- **Data Ingestion:** Continuously collect DNS traffic data from network logs or packet captures.

- **Preprocessing:** Apply the same preprocessing steps used during training to the incoming data.

- **Model Inference:** Use the trained model to classify each DNS query as benign or malicious in real-time.

*Example:*

```
import tensorflow as tf
```

*# Load the trained model*

```
model = tf.keras.models.load_model('path/to/model')
```

```python
# Real-time data ingestion and preprocessing
def preprocess_data(data):
    # Apply preprocessing steps
    return preprocessed_data


def predict(data):
    preprocessed_data = preprocess_data(data)
    predictions = model.predict(preprocessed_data)
    return predictions


# Example of real-time data processing
while True:
    data = get_real_time_dns_data()  # Function to get real-time DNS data
    predictions = predict(data)
    handle_predictions(predictions)  # Function to handle the predictions
```

# Results and Discussions

## 7.1 Model Performance Metrics

1. *Accuracy*

   - **Definition:** Accuracy measures the proportion of correctly classified instances out of the total instances.

   - **Formula:** Accuracy=TP+TNTP+TN+FP+

   - **Example Code:**

   from sklearn.metrics import accuracy_score

   accuracy = accuracy_score(y_test, y_pred_classes)

   print(f"Accuracy: {accuracy:.2f}")

2. *Precision*

   - **Definition:** Precision measures the proportion of true positive instances out of the total instances classified as positive.

   - **Formula:** Precision=TP+FP

3. *Recall*

   - **Definition:** Recall measures the proportion of true positive instances out of the total actual positive instances.

   - **Formula:** Recall=TPTP+FNRecall=$TP+FNTP$

4. *F1-Score*

   - **Definition:** The F1-score is the harmonic mean of precision and recall.

   - **Formula:** F1-Score=2×Precision×RecallPrecision+RecallF1-Score=2×Precision+RecallPrecision×Recall

5. *Confusion Matrix*

   - **Description:** The confusion matrix summarizes the performance of the model by showing the true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN).

## 7.2 Comparative Analysis

1. *Performance Comparison*

   - Present a table or chart comparing the performance metrics (accuracy, precision, recall, F1-score) of different models.

   - **Discussion:**

     - **CNN:** Achieved high accuracy and precision but struggled with recall, indicating it may miss some malicious instances.

     - **RNN:** Performed well across all metrics, particularly in recall, making it effective for detecting sequential anomalies.

- **Autoencoder:** Showed good performance in anomaly detection but had a higher false positive rate, indicating a need for further tuning.

## 7.3 Real-Time Detection Performance

1. *Latency*

   - **Description:** Measure the time taken by the model to process and classify DNS queries in real-time.

   - **Example Code:**

   import time

   start_time = time.time()

   predictions = model.predict(real_time_data)

   end_time = time.time()

   latency = end_time - start_time

   print(f"Latency: {latency:.2f} seconds")

2. *Scalability*

   - **Description:** Test the model's ability to handle large volumes of DNS traffic, ensuring it can scale horizontally.

3. *False Positives/Negatives*

   - **Description:** Analyze the rate of false positives and false negatives in real-time detection.

   - **Example Code:**

   false_positives = sum((y_real_time == 0) & (y_pred_real_time == 1))

   false_negatives = sum((y_real_time == 1) & (y_pred_real_time == 0))

   print(f"False Positives: {false_positives}")

   print(f"False Negatives: {false_negatives}")

# Challenges and Limitations

## 8.1 Data Quality and Representativeness

Addressing the issues related to the quality and representativeness of the data used for training and evaluation.

1. ***Diversity of Data***

   - The dataset may not fully capture the diversity of real-world DNS traffic, potentially limiting the model's ability to generalize to new, unseen data.

     - If the data is collected from a limited number of sources, it may not represent the full spectrum of DNS traffic patterns encountered in different environments.

     - The dataset may be biased towards certain types of traffic or attack patterns, affecting the model's performance on other types.

     - Future work should focus on collecting more diverse and representative datasets from various sources, including different geographic regions, network types, and user behaviors.

2. ***Labeling Accuracy***

   - The accuracy of the labels used for training is critical. Any errors or inconsistencies in labeling can negatively impact the model's performance.

     - Manual labeling by experts can be time-consuming and prone to human error.

     - Automated labeling tools may not always be accurate, especially for new or sophisticated attack patterns.

     - Implement rigorous validation processes for labeling, combining manual and automated methods to ensure high-quality, accurate labels.

## 8.2 Model Complexity and Resource Requirements

Discussing the challenges related to the complexity of deep learning models and the computational resources required for training and inference.

1. ***Computational Resources***

   - Deep learning models can be resource-intensive, requiring significant computational power for training and inference.

     - Training deep learning models can take a considerable amount of time, especially with large datasets and complex architectures.

     - High-performance hardware (e.g., GPUs) is often required to train and deploy deep learning models efficiently.

     - Explore more efficient model architectures and optimization techniques, such as model pruning, quantization, and distributed training, to reduce resource requirements.

2. ***Model Interpretability***

   - The complexity of deep learning models can make them difficult to interpret, potentially limiting their usefulness for security analysts.

- Deep learning models are often considered "black boxes," making it challenging to understand how they arrive at their decisions.

- Providing clear and actionable explanations for the model's predictions is crucial for gaining trust and ensuring effective threat response.

- Develop and integrate explainable AI (XAI) techniques, such as SHAP, LIME, and attention mechanisms, to improve model interpretability and provide insights into the decision-making process.

## 8.3 False Positives and False Negatives

Addressing the challenges related to balancing precision and recall to minimize false positives and false negatives.

1. *False Positives*

   - High false positive rates can overwhelm security analysts with alerts, reducing the effectiveness of the detection system.

     - Excessive false positives can lead to alert fatigue, where analysts become desensitized to alerts and may overlook genuine threats.

     - Investigating false positives consumes valuable time and resources that could be better spent on addressing actual threats.

     - Further tuning and optimization of the model, including adjusting detection thresholds and incorporating additional contextual information, can help reduce false positives.

2. *False Negatives*

   - Missing malicious instances (false negatives) can have serious security implications, as undetected threats can cause significant damage.

     - False negatives represent missed opportunities to detect and mitigate threats, potentially leading to security breaches.

     - Ensuring high recall is essential for effective threat detection, but increasing sensitivity can also lead to more false positives.

     - Implement ensemble methods or hybrid models that combine multiple detection techniques to improve recall without significantly increasing false positives.

## 8.4 Real-Time Detection and Scalability

Discussing the challenges related to implementing the model in real-time detection scenarios and ensuring scalability.

1. *Real-Time Processing*

   - Implementing the model in a real-time processing pipeline requires efficient data ingestion, preprocessing, and inference.

     - Ensuring low latency in processing DNS queries is crucial for timely threat detection and response.

     - Handling high volumes of DNS traffic in real-time can be challenging, especially in large-scale network environments.

- Optimize the model and preprocessing steps for speed, and explore distributed processing frameworks (e.g., Apache Kafka, Apache Spark) to handle high traffic volumes efficiently.

2. ***Scalability***

- Ensuring that the model can scale horizontally to accommodate increasing traffic loads is essential for deployment in real-world environments.

  - Scaling the model requires efficient allocation of computational resources, including load balancing and resource management.

  - Integrating the model with existing network infrastructure and security systems can be complex and require careful planning.

  - Implement scalable architectures and distributed processing techniques, and develop strategies for seamless integration with existing systems.

## 8.5 Adaptability to Evolving Threats

Addressing the challenges related to the model's ability to adapt to new and evolving threats.

1. ***Emerging Threats***

- Cyber threats are constantly evolving, with attackers developing new techniques to bypass detection systems.

  - A static model may become less effective over time as new attack patterns emerge.

  - Ensuring that the model can continuously learn and adapt to new threats is crucial for maintaining its effectiveness.

  - Implement mechanisms for continuous learning and updating the model with new data, including online learning and incremental updates, to ensure it remains effective against emerging threats.

# Conclusion

## 9.1 Key Findings

- **RNN Model:** Demonstrated the highest overall performance, particularly excelling in recall, making it effective for detecting sequential anomalies.

- **CNN Model:** Showed strong precision but had a slightly lower recall, indicating potential misses of some malicious instances.

- **Autoencoder:** Was effective in anomaly detection but had a higher false positive rate, suggesting a need for further tuning.

## 9.2 Challenges and Limitations

- **Data Quality:** The dataset may not fully capture the diversity of real-world DNS traffic, potentially limiting the model's generalizability.

- **Model Complexity:** Deep learning models are resource-intensive and can be difficult to interpret, posing challenges for real-time deployment and usability by security analysts.

- **False Positives/Negatives:** Balancing precision and recall to minimize false positives and false negatives remains a significant challenge.

- **Scalability:** Ensuring the model can handle large volumes of DNS traffic efficiently is crucial for real-world deployment.

- **Adaptability:** Continuous learning and adaptation to new and evolving threats are essential for maintaining the model's effectiveness.

## 9.3 Future Work

- **Hybrid Models:** Combining the strengths of different architectures to improve detection accuracy.

- **Real-Time Detection:** Optimizing the model for low-latency inference and implementing distributed processing frameworks.

- **Continuous Learning:** Developing techniques for online and active learning to adapt to emerging threats.

- **Model Interpretability:** Integrating explainable AI techniques to provide clear and actionable insights.

- **Integration with Security Systems:** Investigating the integration of DNS detection with other security systems for a comprehensive threat detection approach.

## 9.4 Practical Implications

- **Deployment:** Ensuring the necessary infrastructure and seamless integration with existing security systems.

- **Operational Challenges:** Developing strategies for managing alerts and maintaining the model's performance over time.

# Coding

### 10.1 Data Loading and Initial Exploration

```python
import pandas as pd

import numpy as np

import matplotlib.pyplot as plt

import seaborn as sns


# Load the data with custom column names

df = pd.read_csv('data.csv', names=["Timestamp", "Duration", "Type", "Level",
                "Client", "Client ID", "Query ID",
                "Query Name", "View", "Recursion",
                "Query Type", "Query", "Class",
                "Record Type", "Flags", "IP Address"], low_memory=False)


# Display the first few rows of the DataFrame

print("Initial DataFrame:")

print(df.head())


# Get the shape of the DataFrame

data_shape = df.shape


# Display the shape

print(f"Number of rows: {data_shape[0]}")

print(f"Number of columns: {data_shape[1]}")


# Check for null values

null_values = df.isnull().sum()

print("Null values in each column:")

print(null_values)


# Display rows with any null values
```

```python
rows_with_nulls = df[df.isnull().any(axis=1)]
print("Rows with null values:")
print(rows_with_nulls)


# Drop rows with any null values
df.dropna(inplace=True)


# Verify that the rows with null values have been dropped
print("Null values after dropping:")
print(df.isnull().sum())
print("Shape after dropping nulls:", df.shape)


# List of columns to select
columns_to_select = ["Timestamp", "Duration", "Client", "Client ID", "Query ID", "Query", "Class", "Record Type", "Flags", "IP Address"]


# Select the specified columns
selected_df = df[columns_to_select]


# Display the first few rows of the selected DataFrame
print("Selected DataFrame:")
print(selected_df.head())


# Create a copy of the DataFrame
selected_columns = selected_df.copy()
print("Shape of selected columns:", selected_columns.shape)
```

**10.2 Time Formatting**

```python
import dask.dataframe as dd


# Define the function to format time
def format_time(time_str):
    if time_str is None or time_str == ":
```

```python
        return None  # Handle NA values

    try:
        # Check if the time string has hours
        if len(time_str.split(':')) == 3:
            # Split the time string into hours, minutes, seconds, and milliseconds
            hours, minutes, seconds_milliseconds = time_str.split(':')
            seconds, milliseconds = seconds_milliseconds.split('.')
        elif len(time_str.split(':')) == 2:
            # Split the time string into minutes, seconds, and milliseconds
            minutes, seconds_milliseconds = time_str.split(':')
            seconds, milliseconds = seconds_milliseconds.split('.')
            hours = 0
        else:
            print(f"Invalid time format: {time_str}")
            return None

        # Convert to integers
        hours = int(hours)
        minutes = int(minutes)
        seconds = int(seconds)
        milliseconds = int(milliseconds)

        # Calculate total milliseconds
        total_milliseconds = (hours * 60 * 60 * 1000) + (minutes * 60 * 1000) + (seconds * 1000) + milliseconds

        # Convert to formatted time string
        formatted_time = f"{hours:02}:{minutes:02}:{seconds:02}.{milliseconds:03}"
        return formatted_time
    except Exception as e:
        print(f"Error formatting time: {time_str} - {e}")
        return None  # Return None for invalid entries
```

```python
# Apply the function to each row in the 'Duration' column
selected_columns['time'] = selected_columns['Duration'].map(format_time, meta=('time', 'object'))


# Drop the 'Duration' column
selected_columns = selected_columns.drop('Duration', axis=1)


# Print the updated DataFrame (compute to ensure data is loaded)
print("Updated DataFrame with formatted time:")
print(selected_columns.compute())
```

## 10.3 Regex Pattern Matching

```python
import re


# Define the regular expressions
regex_patterns = {
  "service_specific_dns": r"^_ldap\._tcp\.dc\._msdcs\.[A-Z]+\.[A-Z]+$",

  "uuid_based_service_specific_dns": r"^_ldap\._tcp\.[0-9a-fA-F-]{36}\.domains\._msdcs\.[A-Z]+\.[A-Z]+$",

  "standard_hostname_domain": r"^[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+\.[a-zA-Z]{2,}$",

  "complex_subdomain_structure": r"^[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+\.[a-zA-Z]{2,}$",

  "standard_hostname_subdomain": r"^[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+\.[a-zA-Z]{2,}$",

  "simple_domain": r"^[a-zA-Z0-9-]+\.[a-zA-Z]{2,}$",  # Added to match simple domains like "amazon.com"

  "single_label_domain": r"^[a-zA-Z0-9-]+$",  # Added to match single label domains like "."

  "extended_subdomain_structure": r"^[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+\.[a-zA-Z]{2,}$",  # Extended subdomain structure

  "another_extended_subdomain_structure": r"^[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+\.[a-zA-Z]{2,}$",  # Another extended subdomain structure

  "ip_based_subdomain": r"^\d{1,3}\.\d{1,3}\.\d{1,3}\.\d{1,3}\.[a-zA-Z0-9-]+\.[a-zA-Z]{2,}$",  # IP-based subdomain

  "complex_ip_based_subdomain": r"^\d{1,3}\.\d{1,3}\.\d{1,3}\.\d{1,3}\.[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+\.[a-zA-Z]{2,}$",  # Complex IP-based subdomain

  "subdomain_with_underscore": r"^[a-zA-Z0-9-_]+\.[a-zA-Z0-9-]+\.[a-zA-Z0-9-]+\.[a-zA-Z]{2,}$",  # Subdomain with underscore

  "complex_subdomain_with_underscore": r"^[a-zA-Z0-9-_]+\.[a-zA-Z0-9-_]+\.[a-zA-Z0-9-_]+\.[a-zA-Z0-9-_]+\.[a-zA-Z]{2,}$"  # Complex subdomain with underscore
```

```python
}

# Extract query names
query_names = selected_columns['Query'].tolist()


# Function to match query names to regex patterns
def match_query_names(query_names, regex_patterns):
  matches = {}
  for name in query_names:
      for pattern_name, pattern in regex_patterns.items():
          if re.match(pattern, name):
              matches[name] = pattern_name
              break
  return matches


# Match the query names
matched_queries = match_query_names(query_names, regex_patterns)


# Print the results
print("Matched Queries:")
for query, pattern_name in matched_queries.items():
  print(f"Query: {query} matches pattern: {pattern_name}")


# Calculate the percentage of matched queries
total_queries = len(query_names)
matched_count = len(matched_queries)
percentage_matched = (matched_count / total_queries) * 100 if total_queries > 0 else 0


# Print the results
print(f"Total Queries: {total_queries}")
print(f"Matched Queries: {matched_count}")
print(f"Percentage of Queries Matched: {percentage_matched:.2f}%")
```

```python
# Get unique matched queries
matched_set = set(matched_queries.keys())


# Print the unique matched queries
print("Unique Queries that matched a pattern:")
for query in matched_set:
    print(query)


# Get unique unmatched queries
unmatched_queries = set(query_names) - matched_set


# Print the unique unmatched queries
print("Unique Queries that did not match any pattern:")
for query in unmatched_queries:
    print(query)


# Add new regex patterns for the unique unmatched queries
for query in unmatched_queries:
    # Create a regex pattern for each unique unmatched query
    regex_patterns[f"unique_query_{query}"] = re.escape(query)


# Print the updated regex patterns
print("Updated Regex Patterns:")
for name, pattern in regex_patterns.items():
    print(f"{name}: {pattern}")
```

## 10.4 Domain and IP Address Extraction

```python
# Extracting the IP Address and saving it in IP_list.txt
def extract_ip_addresses(ip_string):
    # Use regex to find all valid IPv4 and IPv6 addresses
    ip_pattern = r'(\d{1,3}\.\d{1,3}\.\d{1,3}\.\d{1,3}|[0-9a-fA-F:]+)'
    return re.findall(ip_pattern, ip_string)
```

```python
# Extract IP addresses from the 'IP Address' column
ip_addresses = selected_columns['IP Address'].apply(extract_ip_addresses)


# Flatten the list of lists and get unique IP addresses
unique_ips = set(ip for sublist in ip_addresses for ip in sublist)


# Save the unique IP addresses to a text file
with open('IP_list.txt', 'w') as file:
    for ip in unique_ips:
        file.write(f"{ip}\n")  # Write each IP on a new line


print("Unique IP addresses saved to IP_list.txt")


# Extracting the domain and saving it in domain_list.txt
def extract_domain(query):
    # Check if the query is already a simple domain
    if re.match(r'^[a-zA-Z0-9-]+\.[a-zA-Z]{2,}$', query):
        return query  # Return the query as it is a valid domain


    # Use regex to match the domain part
    match = re.search(r'([a-zA-Z0-9-]+\.[a-zA-Z0-9-]+\.[a-zA-Z]{2,})$', query)
    if match:
        return match.group(1)  # Return the matched domain
    return None  # Return None if no match is found


# Apply the function to the 'Query' column
selected_columns['Domain'] = selected_columns['Query'].apply(extract_domain)


# Check for null values in the 'Domain' column
if selected_columns['Domain'].isnull().any():
    print("The 'Domain' column contains null values.")
else:
    print("The 'Domain' column does not contain null values.")
```

```python
# Get the count of null values in the 'Domain' column
null_count = selected_columns['Domain'].isnull().sum()
print(f"The 'Domain' column contains {null_count} null values.")


# Save the domains to a text file
with open('domain_list.txt', 'w') as file:
  for domain in selected_columns['Domain'].unique():  # Get unique domains
    file.write(f"{domain}\n")  # Write each domain on a new line


print("Domains saved to domain_list.txt")


# Print the updated DataFrame with the extracted domains
print("Updated DataFrame with extracted domains:")
print(selected_columns[['Query', 'Domain']])


# Calculate the percentage of queries that were successfully extracted
num_extracted_queries = selected_columns['Domain'].notnull().sum()
total_queries = len(selected_columns)
percentage_extracted_queries = (num_extracted_queries / total_queries) * 100


print(f"Percentage of queries extracted: {percentage_extracted_queries:.2f}%")
```

**10.5 Merging Domain Data**

```python
# Load the domain labels file
df_labels = pd.read_csv('/content/domain_labels.txt',
            delimiter=',',
            header=None,
            names=['Domain', 'Categorized', 'Category', 'Risk'])


# Display the first few rows of the DataFrame
print("Domain Labels DataFrame:")
print(df_labels.head())
```

```python
# Merge the extracted domains with the existing domain data
merged_df = pd.merge(selected_columns, df_labels, on='Domain', how='inner', indicator=True)
merged_df = merged_df.drop(columns=['_merge'])


# Display the merged DataFrame
print("Merged DataFrame:")
print(merged_df.head())
```

## 10.6 Mapping Categories and Risks

```python
# Define mappings for Category and Risk
category_mapping = {
  'Malicious': 1,
  'Suspicious': 0.5,  # Optional if you want to include this category
  'Benign': 0
}


risk_mapping = {
  'High': 2,
  'Medium': 1,
  'Low': 0
}


# Convert Category and Risk to numerical values
merged_df['Category'] = merged_df['Category'].replace(category_mapping)
merged_df['Risk'] = merged_df['Risk'].replace(risk_mapping)


# Display the updated DataFrame
print("Updated DataFrame with Numerical Categories and Risks:")
print(merged_df[['Domain', 'Category', 'Risk']])
```

## 10.7 Data Cleaning and Resampling

```python
# Clean up the 'IP Address' column by removing parentheses
```

```python
merged_df['IP Address'] = merged_df['IP Address'].str.replace(r'[()]', '', regex=True)

# Set the 'Timestamp' column as the index and ensure it's in datetime format
merged_df['Timestamp'] = pd.to_datetime(merged_df['Timestamp'])
merged_df.set_index('Timestamp', inplace=True)


# Get unique IP addresses in the DataFrame
unique_ip_addresses = merged_df['IP Address'].unique()


# Iterate through each unique IP address
for ip_address in unique_ip_addresses:
  # Filter the DataFrame by the current IP address
  filtered_df = merged_df[merged_df['IP Address'] == ip_address]


  # Ensure the DataFrame is sorted by index (timestamp) before resampling
  filtered_df = filtered_df.sort_index()


  # Resample the data to the desired time intervals and count the number of queries
  resampled_30min = filtered_df.resample('30T').size()
  resampled_1hour = filtered_df.resample('1H').size()
  resampled_1day = filtered_df.resample('1D').size()
  resampled_3days = filtered_df.resample('3D').size()


  # Combine the results into a DataFrame
  results = pd.DataFrame({
      '30min': resampled_30min,
      '1hour': resampled_1hour,
      '1day': resampled_1day,
      '3days': resampled_3days
  }).fillna(0).astype(int)


  # Print the results for the current IP address
  print(f"\nResampled results for IP address {ip_address}:")
```

```
    print(results)
```

**10.8 Data Preprocessing for Machine Learning**

```python
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.model_selection import train_test_split
import ipaddress

# Encode categorical variables into numerical variables
le_client_id = LabelEncoder()
le_query = LabelEncoder()
le_domain = LabelEncoder()
le_class = LabelEncoder()
le_record_type = LabelEncoder()
le_flags = LabelEncoder()

merged_df['Client ID'] = le_client_id.fit_transform(merged_df['Client ID'])
merged_df['Query'] = le_query.fit_transform(merged_df['Query'])
merged_df['Domain'] = le_domain.fit_transform(merged_df['Domain'])
merged_df['Class'] = le_class.fit_transform(merged_df['Class'])
merged_df['Record Type'] = le_record_type.fit_transform(merged_df['Record Type'])
merged_df['Flags'] = le_flags.fit_transform(merged_df['Flags'])

# Drop the 'Client' column
merged_df = merged_df.drop('Client', axis=1)

# Convert IP Address to numerical values
def ip_to_int(ip):
  try:
    return int(ipaddress.IPv4Address(ip))
  except ValueError:
    return int(ipaddress.IPv6Address(ip))

merged_df['IP Address'] = merged_df['IP Address'].apply(ip_to_int)
```

```python
# Convert time to seconds since start of day
def time_to_seconds(time):
  h, m, s = time.split(':')
  if '.' in s:
    s, ms = s.split('.')
  else:
    ms = '0'
  return int(h) * 3600 + int(m) * 60 + int(s) + int(ms) / 1000


merged_df['time'] = merged_df['time'].apply(time_to_seconds)


# Verify there are no object types left (besides the label 'Categorized')
print("Data Types After Encoding:")
print(merged_df.dtypes)


# Scale the numerical data
scaler = StandardScaler()
merged_df[['Query ID', 'IP Address', 'time']] = scaler.fit_transform(merged_df[['Query ID', 'IP Address', 'time']])


# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(merged_df.drop('Categorized', axis=1),
merged_df['Categorized'], test_size=0.2, random_state=42)


# Print the shapes of the training and testing sets
print("Training set shape:", X_train.shape, y_train.shape)
print("Testing set shape:", X_test.shape, y_test.shape)
```

# 10.9: Deep Learning Model Training

```python
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Dropout, Input


# Define the model architecture
```

```python
model = Sequential()
model.add(Input(shape=(X_train.shape[1],)))  # Use Input layer to specify the input shape
model.add(Dense(64, activation='relu'))
model.add(Dropout(0.2))
model.add(Dense(32, activation='relu'))
model.add(Dropout(0.2))
model.add(Dense(1, activation='sigmoid'))


# Compile the model
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])


# Train the model
model.fit(X_train, y_train, epochs=10, batch_size=128, validation_data=(X_test, y_test))


# Evaluate the model
loss, accuracy = model.evaluate(X_test, y_test)
print(f'Test Loss: {loss:.3f}')
print(f'Test Accuracy: {accuracy:.3f}')
```

**10.10 Convolutional Neural Network (CNN) Model Training**

```python
from tensorflow.keras.layers import Conv1D, MaxPooling1D, Flatten


# Reshape data to add an extra dimension for Conv1D
X = merged_df.drop('Categorized', axis=1).values
X = X.reshape((X.shape[0], X.shape[1], 1))  # Adding an extra dimension
y = merged_df['Categorized'].values


# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)


# Define the CNN model architecture
model = Sequential()
model.add(Input(shape=(X_train.shape[1], 1)))  # Use Input layer to specify the input shape
```

```python
model.add(Conv1D(64, kernel_size=2, activation='relu'))
model.add(MaxPooling1D(pool_size=2))
model.add(Dropout(0.2))
model.add(Flatten())
model.add(Dense(32, activation='relu'))
model.add(Dropout(0.2))
model.add(Dense(1, activation='sigmoid'))


# Compile the model
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])


# Train the model
model.fit(X_train, y_train, epochs=10, batch_size=128, validation_data=(X_test, y_test))


# Make predictions on the test set
y_pred = model.predict(X_test)
y_pred = (y_pred > 0.5).astype(int)  # Convert probabilities to binary outputs (0 or 1)


# Calculate and print metrics
accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred)
recall = recall_score(y_test, y_pred)
f1 = f1_score(y_test, y_pred)
conf_matrix = confusion_matrix(y_test, y_pred)
class_report = classification_report(y_test, y_pred)


print(f"CNN Accuracy: {accuracy}")
print(f"CNN Precision: {precision}")
print(f"CNN Recall: {recall}")
print(f"CNN F1-Score: {f1}")
print(f"CNN Confusion Matrix:\n{conf_matrix}")
print(f"CNN Classification Report:\n{class_report}")
```

## 10.11 Long Short-Term Memory (LSTM) Model Training

```python
from tensorflow.keras.layers import LSTM

# Reshape data to add an extra dimension for LSTM
X = merged_df.drop('Categorized', axis=1).values
X = X.reshape((X.shape[0], X.shape[1], 1))  # Adding an extra dimension
y = merged_df['Categorized'].values

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Define the LSTM model architecture
model = Sequential()
model.add(Input(shape=(X_train.shape[1], 1)))  # Input shape for LSTM
model.add(LSTM(64, return_sequences=True))
model.add(Dropout(0.2))
model.add(LSTM(32))
model.add(Dropout(0.2))
model.add(Dense(1, activation='sigmoid'))

# Compile the model
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])

# Train the model
model.fit(X_train, y_train, epochs=10, batch_size=128, validation_data=(X_test, y_test))

# Make predictions on the test set
y_pred = model.predict(X_test)
y_pred = (y_pred > 0.5).astype(int)  # Convert probabilities to binary outputs (0 or 1)

# Calculate and print metrics
accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred)
```

```python
recall = recall_score(y_test, y_pred)
f1 = f1_score(y_test, y_pred)
conf_matrix = confusion_matrix(y_test, y_pred)
class_report = classification_report(y_test, y_pred)


print(f"LSTM Accuracy: {accuracy}")
print(f"LSTM Precision: {precision}")
print(f"LSTM Recall: {recall}")
print(f"LSTM F1-Score: {f1}")
print(f"LSTM Confusion Matrix:\n{conf_matrix}")
print(f"LSTM Classification Report:\n{class_report}")
```

**10.12 Autoencoder Model Training**

```python
from tensorflow.keras.models import Model
from tensorflow.keras.layers import Input, Dense
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
from sklearn.metrics import precision_score, recall_score, f1_score, confusion_matrix
import numpy as np


# Assuming merged_df is your DataFrame and 'Categorized' is the label column
# Drop the label column for unsupervised learning
X = merged_df.drop('Categorized', axis=1).values
y = merged_df['Categorized'].values  # Assuming 'Categorized' is binary: 1 for anomaly, 0 for normal


# Scale the data
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)


# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2, random_state=42)


# Define the autoencoder model architecture
```

```python
input_dim = X_train.shape[1]
encoding_dim = 32  # Size of the encoded representation


# Input layer
input_layer = Input(shape=(input_dim,))


# Encoder layers
encoded = Dense(encoding_dim, activation='relu')(input_layer)
encoded = Dense(16, activation='relu')(encoded)


# Decoder layers
decoded = Dense(encoding_dim, activation='relu')(encoded)
decoded = Dense(input_dim, activation='sigmoid')(decoded)


# Autoencoder model
autoencoder = Model(inputs=input_layer, outputs=decoded)


# Compile the model
autoencoder.compile(optimizer='adam', loss='mean_squared_error')


# Train the model
autoencoder.fit(X_train, X_train,
        epochs=50,
        batch_size=256,
        shuffle=True,
        validation_data=(X_test, X_test))


# Use the autoencoder to predict the test set
X_test_pred = autoencoder.predict(X_test)


# Calculate the reconstruction error
reconstruction_error = np.mean(np.power(X_test - X_test_pred, 2), axis=1)
```

```python
# Determine a threshold for anomaly detection
threshold = np.percentile(reconstruction_error, 95)  # 95th percentile


# Predict anomalies based on the threshold
y_pred = (reconstruction_error > threshold).astype(int)


# Calculate and print metrics
precision = precision_score(y_test, y_pred)
recall = recall_score(y_test, y_pred)
f1 = f1_score(y_test, y_pred)
conf_matrix = confusion_matrix(y_test, y_pred)


print(f"Precision: {precision:.2f}")
print(f"Recall: {recall:.2f}")
print(f"F1-Score: {f1:.2f}")
print(f"Confusion Matrix:\n{conf_matrix}")
```

## 10.13 Website Scraping for Domain Categorization

```python
import requests
from bs4 import BeautifulSoup
import time


class DomainLabel():
    '''
    Class to label domain names using the McAfee SmartFilter Internet / Webwasher URL Filter Database
    '''
    def __init__(self):
```

### 10.13.1 Initialize Browser Simulation

```python
        # =====================
        # Set the configuration of a simulated browser
        self.headers = {'User-Agent' : 'Mozilla/5.0 (Windows; Ryzen Windows 23H2)',
```

```
                'Accept' :
'text/html,application/xhtml+xml,application/xml;q=0.9,image/webp,image/apng,*/*;q=0.8',
                'Accept-Language' : 'en-US,en;q=0.9,de;q=0.8'
            }
    self.base_url = 'http://www.trustedsource.org/sources/index.pl'

    # ==========
```

### 10.13.2 Retrieve Tokens for Form Submission

```
    r = requests.get(self.base_url, headers=self.headers)

    bs = BeautifulSoup(r.content, "html.parser")

    form = bs.find("form", {"class": "contactForm"})

    # ==========


    self.token1 = form.find("input", {'name': 'e'}).get('value')

    self.token2 = form.find("input", {'name': 'c'}).get('value')

    self.headers['Referer'] = self.base_url


 def lookup(self, url):
    # Subsection 10.13.3: Domain Lookup
    payload = {'e': (None, self.token1),
            'c': (None, self.token2),
            'action': (None, 'checksingle'),
            'product': (None, '01-ts'),
            'url': (None, url)}


    r = requests.post('https://www.trustedsource.org/en/feedback/url', headers=self.headers, files=payload)

    bs = BeautifulSoup(r.content, "html.parser")

    table = bs.find("table", {"class": "result-table"})


    # ===================

    # When the database return 'None', sleep for 3 sec and retry (w/ new configurations)

    while table is None:

        print('-Retry-')
```

```
        del r

        time.sleep(3)

        # ==========
```

### 10.13.4 Update Browser Configuration

```
        self.update() # Update the configurations of the simulated browser

        payload = {'e': (None, self.token1),

                'c': (None, self.token2),

                'action': (None, 'checksingle'),

                'product': (None, '01-ts'),

                'url': (None, url)}

        r = requests.post('https://www.trustedsource.org/en/feedback/url', headers=self.headers, files=payload)

        bs = BeautifulSoup(r.content, "html.parser")

        table = bs.find("table", {"class": "result-table"})


    td = table.find_all('td')

    categorized = td[len(td)-3].text # Whether the given domain is categorized in the database

    category = td[len(td) - 2].text[2:] # The category of the given domain

    risk = td[len(td) - 1].text # Risk level


    return categorized, category, risk


 def update(self):

    '''

    Function to update the configurations of the simulated browser

    '''
```

### 10.13.5 Update Tokens

```
    # ====================

    r = requests.get(self.base_url, headers=self.headers)

    bs = BeautifulSoup(r.content, "html.parser")

    form = bs.find("form", {"class": "contactForm"})

    # ==========
```

```
self.token1 = form.find("input", {'name': 'e'}).get('value')

self.token2 = form.find("input", {'name': 'c'}).get('value')

self.headers['Referer'] = self.base_url
```

## 10.14 Passing the Extracted TXT File for Domain Categorization

```
import requests

from requests.adapters import HTTPAdapter

from requests.packages.urllib3.util.retry import Retry

import os

import time

from concurrent.futures import ThreadPoolExecutor, as_completed


# =====================

DATA_PATH = 'D:/Anaconda Python/dns detection/data/sample.txt'  # Path of the input domain list

RES_PATH = 'D:/Anaconda Python/dns detection/res/sample.txt'  # Path to save the domain labels


# =====================

class DomainLabel:
  def __init__(self):
```

### *10.14.1 Initialize Session and Load Existing Data*

```
    self.base_url = 'https://www.trustedsource.org/sources/index.pl'  # Use HTTPS

    self.headers = {'User-Agent': 'Mozilla/5.0'}

    self.session = requests.Session()

    retry = Retry(connect=3, backoff_factor=0.5)

    adapter = HTTPAdapter(max_retries=retry)

    self.session.mount('http://', adapter)

    self.session.mount('https://', adapter)

    self.domain_set = set()  # Initialize domain_set here


    # Load already labeled domains

    if os.path.exists(RES_PATH):

        with open(RES_PATH, 'r') as f_input:
```

```python
        for line in f_input.readlines():

            record = line.strip().split(',')

            domain = record[0]

            self.domain_set.add(domain)  # Populate domain_set


    def lookup(self, domain):
        # 10.14.2 Domain Lookup with Error Handling
        try:
            # Set a timeout for the request
            r = self.session.get(self.base_url, headers=self.headers, params={'md5': domain}, timeout=10,
verify=False)  # Disable SSL verification
            if r.status_code == 200:
                categorized = 'True'
                category = 'Unknown'
                risk = 'Unknown'
                # Parse the HTML content to extract the category and risk
                if 'malicious' in r.text:
                    category = 'Malicious'
                    risk = 'High'
                elif 'suspicious' in r.text:
                    category = 'Suspicious'
                    risk = 'Medium'
                else:
                    category = 'Benign'
                    risk = 'Low'
                return categorized, category, risk
            else:
                return 'False', None, None
        except requests.exceptions.Timeout:
            print(f"Timeout error for domain {domain}. Skipping...")
            return 'False', None, None
        except requests.exceptions.RequestException as e:
            print(f"Error while looking up domain {domain}: {e}")
```

```python
        return 'False', None, None


    def save_result(self, domain, categorized, category, risk):
        # 10.14.3 Save Results
        with open(RES_PATH, 'a+') as f_output:
            if categorized == 'False':
                f_output.write('%s,%s\n' % (domain, categorized))
            else:
                f_output.write('%s,%s,%s,%s\n' % (domain, categorized, category, risk))


# =====================
# 10.14.4 Process Domain List
domain_label = DomainLabel()
label_cnt = 0  # Counter of labelled domains


# Read the input domain list
with open(DATA_PATH, 'r') as f_input:
 for line in f_input.readlines():
    domain = line.strip()
    if domain in domain_label.domain_set:  # Check against the instance variable
        label_cnt += 1
        continue


    # Attempt to categorize the domain
    categorized, category, risk = domain_label.lookup(domain)
    print('-Record-#%d %s %s %s %s' % (label_cnt, domain, categorized, category, risk))


    # Save the result
    domain_label.save_result(domain, categorized, category, risk)


    label_cnt += 1
    time.sleep(1)  # Sleep for 1 sec to avoid hitting the API too quickly
```

# References

1. **Published Date**: 2018
   **Title**: "Deep Learning for Malicious Flow Detection"
   **Authors**: Y. Yu, X. Yang, C. Fung
   **Published In**: IEEE International Conference on Communications (ICC)
   **Link**: [IEEE Xplore](IEEE Xplore)

2. **Published Date**: 2019
   **Title**: "DNS Deep Learning: Detecting DGA-based Malware Domains"
   **Authors**: S. Schüppen, M. Wolsing, J. H. Ziegeldorf, K. Wehrle, M. Henze
   **Published In**: IEEE International Conference on Network and Service Management (CNSM)
   **Link**: [IEEE Xplore](IEEE Xplore)

3. **Published Date**: 2020
   **Title**: "DeepDetect: Detection of Malicious Domain Names Based on Deep Learning"
   **Authors**: M. S. Alzahrani, A. A. Ghorbani
   **Published In**: Journal of Network and Computer Applications
   **Link**: [ScienceDirect](ScienceDirect)

4. **Published Date**: 2018
   **Title**: "Detecting Malicious Domain Names Using Deep Learning Approaches"
   **Authors**: S. Lisoněk, M. Rehák
   **Published In**: IEEE Security and Privacy Workshops (SPW)
   **Link**: [IEEE Xplore](IEEE Xplore)

5. **Published Date**: 2019
   **Title**: "Deep Learning for Detecting Malicious Domain Names"
   **Authors**: M. S. Alzahrani, A. A. Ghorbani
   **Published In**: IEEE Access
   **Link**: [IEEE Xplore](IEEE Xplore)

6. **Published Date**: 2020
   **Title**: "A Deep Learning Approach for DNS Tunneling Detection"
   **Authors**: A. Z. Alrawashdeh, C. Purdy
   **Published In**: IEEE Transactions on Network and Service Management
   **Link**: [IEEE Xplore](IEEE Xplore)

7. **Published Date**: 2017
   **Title**: "Detecting Malicious Domains via Graph Inference and Deep Learning"
   **Authors**: Y. Liu, Y. Xiao, H. Zhang, S. Zhang
   **Published In**: IEEE International Conference on Data Mining (ICDM)
   **Link**: [IEEE Xplore](IEEE Xplore)

8. **Published Date**: 2019
   **Title**: "Deep Learning for Detecting Malicious URLs"
   **Authors**: S. Saxe, D. Berlin
   **Published In**: IEEE Security and Privacy Workshops (SPW)
   **Link**: [IEEE Xplore](IEEE Xplore)

9. **Published Date**: 2021
   **Title**: "A Survey on Deep Learning for Cyber Security"

**Authors**: M. A. Ferrag, L. Shu, X. Yang, A. Derhab, L. Maglaras
**Published In**: IEEE Communications Surveys & Tutorials
**Link**: IEEE Xplore

10. **Published Date**: 2020
    **Title**: "Deep Learning for Detecting Malicious Network Traffic"
    **Authors**: J. Zhang, Y. Zhang, X. Chen
    **Published In**: IEEE Transactions on Network and Service Management
    **Link**: IEEE Xplore

# List of figure

# Abbreviations & Acronyms

- AI: Artificial Intelligence

- CNN: Convolutional Neural Network

- DGA: Domain Generation Algorithm

- DNS: Domain Name System

- FP: False Positive

- FN: False Negative

- GPU: Graphics Processing Unit

- GRU: Gated Recurrent Unit

- IDS: Intrusion Detection System

- ISP: Internet Service Provider

- LIME: Local Interpretable Model-agnostic Explanations

- LSTM: Long Short-Term Memory

- MX: Mail Exchange (DNS record type)

- RNN: Recurrent Neural Network

- ROC: Receiver Operating Characteristic

- SHAP: SHapley Additive exPlanations

- SIEM: Security Information and Event Management

- SMOTE: Synthetic Minority Over-sampling Technique

- TP: True Positive

- TN: True Negative

- TXT: Text (DNS record type)

- XAI: Explainable Artificial Intelligence