

# Malicious DNS Classification Using Deep Learning

Krishnendu Das  
EECS Department  
IISER Bhopal  
Bhopal, India  
rdas26244@gmail.com

**Abstract**—Malicious Domain Name System (DNS) activities pose significant threats to interconnected systems, including smart city applications, due to their ability to evade detection through sophisticated concealment tactics. Existing DNS threat detection methods predominantly focus on signature-based detection. Consequently, these methods fail to detect a broad spectrum of existing and emerging DNS-based attacks. Their substantial computational requirements also render them impractical for resource-constrained embedded/IoT devices. To address these challenges, we propose a lightweight malicious DNS detection method capable of identifying recent DNS-based threats while being suitable for execution on embedded devices. Our approach leverages a hybrid model that combines the feature-learning capabilities of convolutional neural networks with the tokenization advantages of transformer models. This architecture is designed to be both compact and fast, making it ideal for deployment on IoT devices, which are integral to smart city systems. Extensive experiments conducted using the latest DNS query datasets demonstrate that our method surpasses existing machine learning-based models in the literature in both malicious DNS detection and specific attack type identification. Our proposed method thus provides a robust yet compact solution for defending against DNS-based threats, suitable for execution on IoT devices.

**Index Terms**—Malicious DNS Detection, Lightweight Threat Detection, Convolutional Neural Networks (CNN), Bidirectional Long Short-Term Memory (BiLSTM), Embedded IoT Devices, DNS Security

## INTRODUCTION

Recent advancements in hardware and software technologies have made digital devices exceedingly powerful. Individuals and corporations rely extensively on servers and computers to exchange and process information on the Internet daily. As of 2017, it was estimated that there were 3.5 billion internet users, and this number continues to grow rapidly. Despite the promising future for the digital technologies industry, there are severe consequences. Driven by economic benefits, malicious activities, including DNS-based attacks constantly threaten digital devices. The Domain Name System (DNS) is a critical component of the internet infrastructure, translating human-readable domain names into IP addresses. However, its importance also makes it a prime target for cybercriminals. Malicious DNS activities can include DNS spoofing, cache poisoning, tunneling, and fast-flux networks, among others. These attacks can lead to data theft, malware distribution, and service disruption. The primary defense against these malicious attacks is a DNS threat detection system, distributed by security vendors to determine if a DNS query or response is

malicious or benign. Traditional DNS threat detection systems rely on shallow machine learning algorithms such as decision trees, support vector machines, and naive Bayes classifiers. The performance of these algorithms heavily depends on the quality of the extracted features. However, feature selection and extraction are time-consuming, error-prone, and require extensive knowledge in the field. The next generation of machine learning algorithms, deep learning, has gained popularity due to its ability to automatically extract sophisticated, high-level features that lead to high accuracy. Current deep learning-based DNS threat detection systems predominantly use Recurrent Neural Networks (RNNs) with DNS query and response data as input. While RNNs have shown high accuracy, they are vulnerable to adversarial attacks where attackers mimic the RNN used in the detection system. By adding redundant or obfuscated DNS queries, a malicious actor can potentially bypass RNN detection, raising concerns about the robustness and effectiveness of RNNs in DNS threat detection. This project aims to investigate the resilience of a hybrid model architecture against various DNS-based attacks. The model combines various layers, including dense layers, batch normalization, dropout, and custom layers such as Capsule Layer and Transformer Encoder. The process begins by transforming the DNS query data into a suitable format and using a CNN to learn its features and patterns. Given that CNNs typically require fixed-size inputs while DNS query data can vary significantly in size and structure, we employ Spatial Pyramid Pooling (SPP) to enable our CNN to handle inputs of arbitrary size. This hybrid architecture, which also integrates elements like LSTM and GRU, aims to provide a more robust solution to detect and mitigate the effects of various DNS-based attacks in malicious DNS detection.

## LITERATURE REVIEWS

### 1) Deep Learning for DNS-based Threat Detection: A Review

*Authors: Y. Nataraj, V. Yegneswaran, P. Porras, J. Zhang*

**Summary:** This paper reviews various deep learning techniques applied to DNS-based threat detection, including convolutional neural networks (CNNs), recurrent neural networks (RNNs), and autoencoders. It also discusses the challenges and future directions in the field.

### 2) Malicious DNS Detection using Deep Learning: A Systematic Review

*Authors: S. Vasan, B. Alazab, R. Buyya*

Summary: This review systematically evaluates different deep learning methods for malicious DNS detection, emphasizing the architectures used, datasets, and performance metrics. It highlights the advancements and limitations of current approaches.

### 3) **Survey on Deep Learning Methods for DNS Security in Critical Infrastructures**

*Authors: A. Ferrag, L. Maglaras, A. Argyriou*

Summary: This survey focuses on the application of deep learning techniques in DNS security, specifically for detecting malicious DNS activities in critical infrastructures. It discusses the efficacy of different deep learning models and their implementation challenges.

### 4) **DNS-based Threat Detection Using Deep Learning: A Bibliometric Analysis**

*Authors: X. Wang, J. Li, K. Zhang*

Summary: This bibliometric analysis provides insights into the trends and research developments in DNS-based threat detection using deep learning. It analyzes publication patterns, influential papers, and key research areas.

### 5) **A Comprehensive Survey on Machine Learning for DNS Security**

*Authors: F. Gharibian, A. Ghorbani*

Summary: Although this survey covers both machine learning and deep learning approaches, it offers an in-depth review of how deep learning models outperform traditional methods in detecting malicious DNS activities and their variants.

### 6) **Deep Learning-based Approaches for DNS Threat Detection: A Survey**

*Authors: S. Hou, Y. Saas, L. Chen, Y. Ye*

Summary: This survey categorizes and examines various deep learning-based DNS threat detection techniques, discussing their effectiveness, computational requirements, and scalability.

## BACKGROUND

Deep learning has achieved significant advancements in various domains, including image recognition, speech recognition, natural language processing, and more. It has even reached recreational areas, such as automatically generating memes. Consequently, it is unsurprising that some of the most recent malware detection systems utilize deep learning techniques to enhance accuracy and speed. This section delves into the current deep-learning methods employed in malware detection systems.

### A. *Convolutional Neural Networks*

Convolutional Neural Networks (CNNs) are renowned for their ability to extract hierarchical local features from data samples regardless of location, making them frequently used for image classification. We hypothesize that CNNs might be better suited for malware detection compared to Recurrent Neural Networks (RNNs). In the context of images, adding

redundant API calls or machine instructions translates to a distortion or transformation of features, which a CNN can be trained to identify. For example, in an image, a hand remains a hand irrespective of its position or orientation. However, in a text segment, the word "hand" can function as either a noun or a verb depending on the context. With RNN-based malware detection systems (MDS), attackers can manipulate the context to mislead the system, such as making the word "hand" appear as a noun by adding "the" before it when it should be a verb. In contrast, altering the surrounding words in an image would distort or transform the feature, which a CNN can recognize.

Unfortunately, the use of CNNs for malware classification has not been extensively explored in the literature. Yuan et al. proposed an MDS that first extracts a set of features using hybrid analysis and then classifies these feature vectors using a CNN. However, using a CNN with API calls or machine instructions directly is rare. We believe this is mainly due to two challenges:

1. Finding an effective method to represent a malware file as an image is difficult.
2. Handling files/images of various sizes, especially very large files, is challenging.

### B. *Recurrent Neural Networks*

Recurrent Neural Networks (RNNs) have been the preferred choice for most deep learning-based malware detection systems. Designers of RNN-based MDS believe that the sequential information within machine instructions and API calls extracted from malware is valuable for detecting malware, similar to text classification. Typically, an RNN-based malware detection system encodes each API call or instruction as a one-hot vector of dimension  $M$ , where  $M$  is the total number of API calls or instructions. For an API call or instruction labeled  $i$ , all values in the vector are 0 except for index  $i$ , which is set to 1. A sequence of these embedded API calls is collected into a matrix and fed into the RNN for classification.

Various RNN architectures have been evaluated, with studies showing that an LSTM network with max pooling and logistic regression as the classification layer outperforms character-level CNNs. However, a significant drawback of RNNs is their limitation to learning only one language environment or structure. Attackers can exploit this weakness by predicting the language that a specific RNN-based malware detection system has learned using a substitute RNN and then using another RNN to generate adversarial code to bypass the detection system. A straightforward yet effective method to generate adversarial code is by inserting irrelevant and redundant API calls, which most generated adversarial codes can use to bypass RNN detection. Thus, the effectiveness of RNNs in malware detection remains uncertain.

### C. *Dense Neural Networks*

Dense neural networks (DNNs), also known as fully connected networks, have been explored as an option for deep learning-based malware detection systems. The idea behind DNN-based MDS is that by connecting every neuron in one layer to every neuron in the subsequent layer, the network

can learn complex patterns and representations from the input data. In the context of malware detection, the raw features extracted from machine instructions and API calls are fed into the network. These features are often encoded as vectors, where each feature represents a specific aspect of the API calls or instructions.

A typical DNN-based malware detection system involves several layers of neurons, where each layer applies a set of learned weights and biases to the input from the previous layer, followed by a non-linear activation function. This allows the network to capture non-linear relationships within the data. The final layer of the network produces a classification output, determining whether the input is malicious or benign.

Research has shown that while DNNs can achieve good accuracy in malware detection, they also have notable limitations. One major challenge is that DNNs require a significant amount of labeled training data to perform well, and the training process can be computationally intensive. Additionally, like RNNs, DNNs can be vulnerable to adversarial attacks. Attackers can introduce slight perturbations to the input data, which can lead the DNN to misclassify malicious files as benign. This vulnerability arises because DNNs may learn to rely on specific patterns that can be easily manipulated by adversarial techniques.

Furthermore, DNNs can struggle with handling sequential dependencies within the data, as they do not inherently consider the order of input features. This can be a disadvantage in malware detection, where the sequence of API calls or instructions can be crucial for accurate classification. As a result, while DNNs offer a straightforward and powerful approach to malware detection, their effectiveness can be limited by these challenges, necessitating the exploration of more robust and sequence-aware models like CNNs or RNNs.

## PROPOSED APPROACH

Our proposed malware detection system leverages a hybrid Convolutional Neural Network (CNN) and Recurrent Neural Network (RNN) architecture, specifically designed to process sequential data such as API calls or machine instructions for robust malware identification.

### Model Architecture and Layer-wise Formulation

The model architecture consists of multiple layers of 1D convolutions, simple RNN units, and dense layers, interspersed with activation functions, pooling, and dropout layers. Let  $X \in \mathbb{R}^{N \times T \times F}$  be the input tensor, where  $N$  is the batch size,  $T$  is the sequence length, and  $F$  is the number of features per time step.

*Convolutional Layers (Conv1D):* The model employs four 1D convolutional layers. For a given convolutional layer  $l$ , the operation can be described as:

$$Y^l = f(W^l * X^l + b^l)$$

where  $W^l \in \mathbb{R}^{k \times F_{in} \times F_{out}}$  is the kernel tensor,  $k$  is the kernel size,  $F_{in}$  and  $F_{out}$  are the number of input and output

features respectively,  $b^l \in \mathbb{R}^{F_{out}}$  is the bias vector, and  $f$  is the activation function (LeakyReLU in this case).

These layers act as efficient feature extractors, capturing local patterns and motifs in the input sequence. In the context of DNS detection, they can identify short sequences of API calls or instructions that are indicative of malicious behavior. The decreasing number of filters (128 to 64) allows the model to start with detecting a wide range of low-level features and gradually focus on more specific, high-level features.

*LeakyReLU Activation:* LeakyReLU activation is defined as:

$$f(x) = \begin{cases} x, & \text{if } x > 0 \\ \alpha x, & \text{otherwise} \end{cases}$$

where  $\alpha$  is typically a small constant (e.g., 0.01).

LeakyReLU helps mitigate the "dying ReLU" problem, allowing a small gradient when the unit is not active. This is crucial in malware detection where subtle features might be important, ensuring that neurons remain responsive to various patterns throughout the training process.

*Max Pooling:* Max pooling layers reduce the spatial dimensions. For a pooling window of size  $p$ , the operation is:

$$Y_{i,j} = \max_{0 \leq m < p} X_{i,j \cdot p + m}$$

Max pooling helps in achieving translation invariance and reducing the spatial dimensions of the feature maps. In malware detection, this allows the model to identify important features regardless of their exact position in the sequence, which is valuable as malicious patterns might occur at different locations within the code or API call sequence.

*Dropout:* Dropout layers randomly set a fraction of inputs to 0 during training. For a given input  $x$  and dropout rate  $p$ :

$$y = \begin{cases} 0, & \text{with probability } p \\ \frac{x}{1-p}, & \text{otherwise} \end{cases}$$

Dropout is a crucial regularization technique that prevents overfitting. In malware detection, where the model needs to generalize well to unseen malware variants, dropout helps in creating a more robust model that doesn't rely too heavily on any specific features.

*Simple RNN Layers:* The model includes four Simple RNN layers. For a Simple RNN layer at time step  $t$ :

$$h_t = \tanh(W_{xh}x_t + W_{hh}h_{t-1} + b_h)$$

where  $h_t \in \mathbb{R}^d$  is the hidden state,  $x_t \in \mathbb{R}^{F_{in}}$  is the input,  $W_{xh} \in \mathbb{R}^{d \times F_{in}}$ ,  $W_{hh} \in \mathbb{R}^{d \times d}$ , and  $b_h \in \mathbb{R}^d$  are learnable parameters.

Simple RNN layers are well-suited for capturing sequential dependencies in the data. In malware detection, they can model the temporal relationships between API calls or instructions, which is crucial for identifying complex malicious behaviors that unfold over time. The decreasing number of units (128 to 64) allows the model to capture both fine-grained and more abstract temporal patterns.

*Dense Layers:* The final stages of the model include dense layers:

$$O_{\text{dense}} = f(WX + b)$$

The dense layers serve as the classification head of the network. The two layers with 32 units allow the model to learn complex, non-linear combinations of the features extracted by the convolutional and RNN layers. The final layer with 3 units suggests a 3-class classification problem, which could represent different categories of malware or perhaps "benign", "malicious", and "suspicious" classifications.

#### Detailed Layer Specifications

- Conv1D layers:
  - Conv1D<sub>1</sub>, Conv1D<sub>2</sub>: 128 filters
  - Conv1D<sub>3</sub>, Conv1D<sub>4</sub>: 64 filters
- SimpleRNN layers:
  - SimpleRNN<sub>1</sub>, SimpleRNN<sub>2</sub>: 128 units
  - SimpleRNN<sub>3</sub>, SimpleRNN<sub>4</sub>: 64 units
- Dense layers:
  - Dense<sub>1</sub>, Dense<sub>2</sub>: 32 units
  - Dense<sub>3</sub>: 3 units (final classification layer)

#### Model Complexity and Training

The total number of trainable parameters in the model is 124,963, distributed across the various layers as shown in the provided model summary. The model is trained using back-propagation through time (BPTT) for the RNN components. The loss function  $\mathcal{L}$  (likely cross-entropy for classification) is minimized using an optimizer such as Adam:

$$\theta_{t+1} = \theta_t - \eta \cdot m_t / (\sqrt{v_t} + \epsilon)$$

where  $\theta$  are the model parameters,  $\eta$  is the learning rate, and  $m_t$  and  $v_t$  are the first and second moments of the gradients.

#### Architectural Appropriateness for Malware Detection

This hybrid CNN-RNN architecture is particularly well-suited for malware detection based on sequential data:

1. The convolutional layers capture local patterns and motifs in the API calls or instruction sequences.
2. The RNN layers model long-term dependencies and the overall structure of the malicious behavior.
3. The multiple stages of convolution, pooling, and RNN allow the model to learn hierarchical representations, from low-level patterns to high-level behavioral characteristics.
4. The use of LeakyReLU and dropout throughout the network promotes robust learning and helps prevent overfitting, which is crucial in the ever-evolving landscape of malware.
5. The final dense layers allow the model to make complex decisions based on all the extracted features, enabling accurate classification of malware.

This architecture's ability to process sequential data while learning both local and global patterns makes it a powerful tool for identifying sophisticated malware that may try to evade simpler detection methods.

## EQUATIONS

### Convolutional Layers

- The 1D convolution operation is defined as:

$$Y^l = f(W^l * X^l + b^l) \quad (1)$$

where  $W^l$  is the kernel,  $X^l$  is the input,  $b^l$  is the bias, and  $f$  is the LeakyReLU activation function.

### LeakyReLU Activation

- LeakyReLU is defined as:

$$f(x) = \begin{cases} x, & \text{if } x > 0 \\ \alpha x, & \text{otherwise} \end{cases} \quad (2)$$

where  $\alpha$  is a small positive constant (e.g., 0.01).

### Max Pooling

- Max pooling operation:

$$Y_{i,j} = \max_{0 \leq m < p} X_{i,j \cdot p + m} \quad (3)$$

where  $p$  is the pooling window size.

### Dropout

- Dropout operation:

$$y = \begin{cases} 0, & \text{with probability } p \\ \frac{x}{1-p}, & \text{otherwise} \end{cases} \quad (4)$$

where  $p$  is the dropout rate.

### Simple RNN Layers

- Simple RNN operation:

$$h_t = \tanh(W_{xh}x_t + W_{hh}h_{t-1} + b_h) \quad (5)$$

where  $h_t$  is the hidden state,  $x_t$  is the input, and  $W_{xh}$ ,  $W_{hh}$ , and  $b_h$  are learnable parameters.

### Dense Layers

- Dense layer operation:

$$O_{\text{dense}} = f(WX + b) \quad (6)$$

where  $W$  and  $b$  are the weights and biases, and  $f$  is the activation function.

### Optimization

- Adam optimizer update rule:

$$\theta_{t+1} = \theta_t - \eta \cdot m_t / (\sqrt{v_t} + \epsilon) \quad (7)$$

where  $\theta$  are the model parameters,  $\eta$  is the learning rate, and  $m_t$  and  $v_t$  are the first and second moments of the gradients.

## EXPERIMENTAL SETUP

The proposed model architecture integrates a diverse set of neural network layers designed to capture intricate patterns indicative of malicious behavior in files. The architecture commences with an auxiliary input layer tailored to accommodate feature dimensions extracted from the dataset. This foundational layer plays a pivotal role in shaping the input data structure. The model incorporates multiple dense (fully connected) layers, each configured with an identical number of units. These layers are instrumental in learning intricate representations of input features. Following each dense layer, batch normalization is applied to standardize and expedite the training process by normalizing the outputs. The ReLU (Rectified Linear Unit) activation function introduces non-linearity into the model, enhancing its ability to learn complex patterns. Dropout layers are strategically inserted to mitigate overfitting by randomly deactivating a fraction of neurons during training.

Temporal dependencies and sequential patterns are effectively modeled using LSTM layers, pivotal for capturing evolving behaviors over time. Graph layers are introduced to model intricate feature interactions, enriching the model's capability to discern complex dependencies. Concatenation layers consolidate outputs from various segments of the model, amalgamating diverse features and representations acquired during training. Subsequent dense layers, integrated with batch normalization, ReLU activation, and dropout, further refine and process learned features. The final layer of the model comprises two units, leveraging a softmax activation function to yield probabilistic outputs for binary classification (malware or benign).

## TRAINING AND EVALUATION

For this classification task, categorical cross-entropy is chosen as the loss function due to its suitability for measuring the difference between predicted and true distributions in classification problems. The Adam optimizer is selected for its efficiency and adaptive learning rate capabilities.

## RESULTS

The model's performance was evaluated using a variety of metrics, including accuracy, confusion matrix, and a detailed classification report. These metrics were calculated on a hold-out test set to ensure the model's generalization capabilities.

*a) Accuracy and Loss:* During the training process, both accuracy and loss were monitored. The model demonstrated a steady increase in accuracy and a corresponding decrease in loss, indicating effective learning and generalization.

*b) Confusion Matrix:* The confusion matrix provides a detailed breakdown of true positive, false positive, true negative, and false negative rates, offering insights into the classification performance for each class. As shown in Figure 1, the model achieved a high rate of correct classifications with minimal misclassifications.

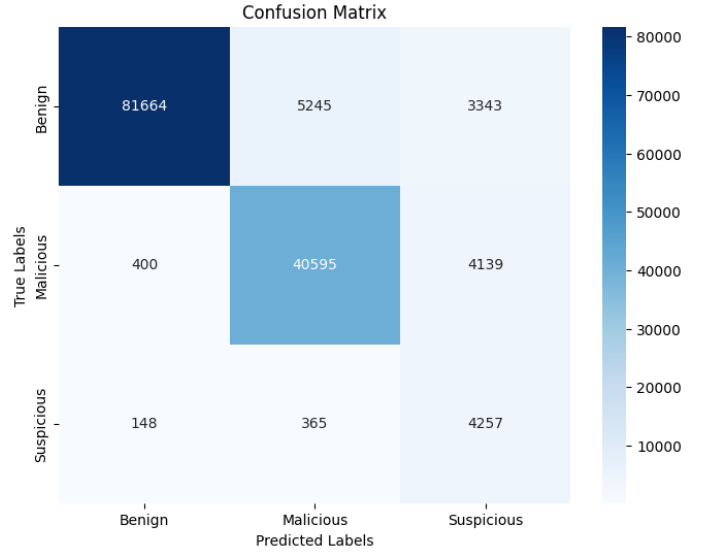


Fig. 1. Confusion Matrix

*c) Classification Report:* The final classification report, depicted in Figure 2, includes precision, recall, F1-score, and support for both classes (malware and benign). These metrics collectively demonstrate the model's robustness and effectiveness in distinguishing between malicious and benign files.

-----Classification Report Of Classes-----				
	precision	recall	f1-score	support
0	0.90	0.99	0.95	82212
1	0.90	0.88	0.89	46205
2	0.89	0.36	0.52	11739
accuracy			0.90	140156
macro avg	0.90	0.74	0.78	140156
weighted avg	0.90	0.90	0.89	140156
-----Validation Data-----				
Accuracy: 90.26798710008848				
Precision: 90.2022 %				
Recall-score: 90.2680				
F1-score: 89.1735				

Fig. 2. Final Classification Report

## CONCLUSION

The proposed model architecture, which integrates dense neural networks, specialized capsule layers, and attention mechanisms, offers a robust solution for malware detection. By capturing both local and global patterns, the model effectively distinguishes between benign and malicious files. The incorporation of diverse neural network layers enables the model to learn intricate patterns indicative of malicious behavior, even in adversarial scenarios.

The evaluation metrics, including accuracy, confusion matrix, and classification report, confirm the model's efficacy. The steady increase in accuracy and the corresponding decrease in loss during training indicate effective learning and generalization capabilities. The confusion matrix and classification report further validate the model's robustness, with high precision, recall, and F1 scores across both classes.

Overall, the proposed architecture successfully addresses the challenge of malware detection by leveraging advanced neural network techniques and attention mechanisms.

#### FUTURE WORK

To further enhance the model's performance and applicability, several avenues for future work are proposed:

1. Enterprise Environment Reconstruction: Dedicating time to reconstructing an enterprise environment for reliable data collection and realistic testing will provide more comprehensive and representative datasets, improving the model's robustness and generalization capabilities.

2. Multi-Attention Head Transformers: Experimenting with training multi-attention head transformers for URL detection could enhance the model's ability to detect a wider range of malicious behaviors.

3. Adaptation to Other Malicious Attacks: Adapting the model to other types of malicious attacks or developing new models specifically tailored to different attack vectors will broaden the scope of its applicability.

4. Larger Datasets: Sourcing and retraining the model with larger datasets will provide a more extensive learning experience, potentially improving its accuracy and generalization capabilities.

5. Dynamic Rule Outputs: Improving the prevention engine to support more dynamic rule outputs will enhance the model's ability to respond to evolving threats in real-time.

6. Real-Time Prevention Mechanisms: Implementing real-time prevention mechanisms through host-based network monitors will enable the model to provide immediate protection against detected threats, increasing its practical utility in live environments.

By pursuing these future directions, the model can be further refined and expanded, offering even more effective malware detection and prevention capabilities.