

# AI Shield: A Hybrid Machine Learning and Deep Learning Approach for Detecting Malicious URLs

CHAITHRA S G<sup>1</sup>, BHAGYASHREE<sup>2</sup>, CHAITHANYA LOKESH<sup>3</sup>, BHARATH. G V<sup>4</sup>, IRFAN KHAN<sup>5</sup>  
<sup>1, 2, 3, 4</sup>*Department of Computer Science and Engineering, Ghousia College of Engineering, Ramanagar, Karnataka, India.*

<sup>5</sup>*Asst. prof, Department of Computer Science and Engineering, Ghousia College of Engineering, Ramanagar, Karnataka, India.*

**Abstract-** *In the modern digital era, the exponential growth of online activities has resulted in an alarming increase in cyber threats, especially phishing and malicious URLs. These threats exploit user trust and vulnerabilities in online systems to steal sensitive credentials and financial information. Traditional defense mechanisms such as blacklist-based filters and static rule-based systems are insufficient, as attackers continuously evolve their techniques to bypass detection. To overcome these limitations, this paper presents AI Shield, a hybrid detection framework that integrates Machine Learning (Decision Tree) and Deep Learning (LSTM) models for the intelligent classification of URLs as safe or malicious. The Decision Tree model performs rule-based lexical analysis, while the LSTM captures sequential dependencies in URL structures, enabling deeper behavioral understanding. Experimental evaluations on the PhishTank 2024 dataset demonstrate that the proposed hybrid approach achieves a detection accuracy of 95%, surpassing standalone ML and DL models. The hybridization approach enhances adaptability, scalability, and real-time detection capability, making AI Shield a robust solution for phishing mitigation in modern cybersecurity infrastructures.*

**Keywords:** *Phishing Detection, Malicious URL, Hybrid Model, Decision Tree, LSTM, Cybersecurity, Deep Learning, Machine Learning.*

## I. INTRODUCTION

The widespread use of the Internet and digital services has made information sharing faster and more convenient. However, this convenience has also led to an increase in cyberattacks, particularly phishing, a

social engineering technique where attackers deceive users into revealing confidential information such as login credentials, banking details, and personal data. These attacks are often executed through malicious URLs, which appear legitimate but redirect users to fraudulent websites. According to the Anti-Phishing Working Group (APWG), phishing attacks have increased by over 60% in recent years, making automated detection systems a necessity.

Traditional approaches to phishing detection, such as blacklists and signature-based filtering, are static and reactive—they depend on previously reported malicious URLs. This limitation makes them ineffective against zero-day attacks, where newly created malicious URLs go undetected. To address this, Machine Learning (ML) and Deep Learning (DL) techniques have been applied to automate phishing detection through pattern recognition and predictive analysis. ML models like Decision Tree, Random Forest, and SVM learn from handcrafted features, while DL models such as CNNs and LSTMs extract deep contextual relationships directly from data.

However, both methods have limitations. ML models often struggle to generalize when faced with unseen data, and DL models require large computational resources and may overfit smaller datasets. To overcome these challenges, this paper proposes AI Shield, a hybrid ML-DL framework that combines the interpretability and simplicity of a Decision Tree with the sequential learning power of an LSTM network. This hybrid approach effectively captures both explicit lexical patterns and hidden contextual dependencies within URLs. By doing so, it enhances detection accuracy, robustness, and adaptability against evolving phishing techniques.

The major contributions of this paper are:

1. Development of a hybrid Decision Tree–LSTM model for malicious URL detection.
2. Extraction of rich lexical and sequential features from URL structures.
3. Experimental validation using the PhishTank 2024 dataset, achieving superior results compared to individual models.
4. Demonstration of the system’s real-time applicability for web browsers and email filtering systems.

## II. LITERATURE REVIEW

The rapid advancement of phishing and malicious URL detection methods has led to extensive research combining traditional and modern computational techniques. Several studies have investigated the application of Machine Learning and Deep Learning models to improve detection efficiency, scalability, and real-time adaptability.

Verma and Singh (2021) implemented Random Forest and SVM classifiers using lexical URL features and achieved high detection accuracy but struggled with false positives. Zhang et al. (2020) applied Convolutional Neural Networks (CNN) for URL text analysis, enabling automatic feature learning but at the cost of heavy computational resources. Patel and Kumar (2022) proposed a hybrid ML-DL approach combining Decision Tree and LSTM models for email phishing detection, which inspired hybrid model research in related areas.

In another study, Almomani et al. (2021) utilised a Bayesianbased phishing detection technique with real-time data analysis, achieving good results for small-scale datasets. Banu et al. (2022) employed ensemble models combining Naïve Bayes and Random Forest for phishing detection, showing that ensemble techniques outperform single classifiers. Similarly, Li and Wang (2023) proposed an attention-based LSTM model that captures semantic relationships within URLs, improving detection precision but requiring significant training data.

Hybrid models have recently emerged as a strong alternative for phishing detection. Khan et al. (2023) developed a hybrid CNNLSTM model for spam

filtering that demonstrated the advantages of combining both sequence-based and convolutional feature extraction. Inspired by these works, AI Shield adopts a Decision Tree–LSTM hybridisation specifically tailored for URL-level phishing detection, where lexical rules and sequential dependencies coexist.

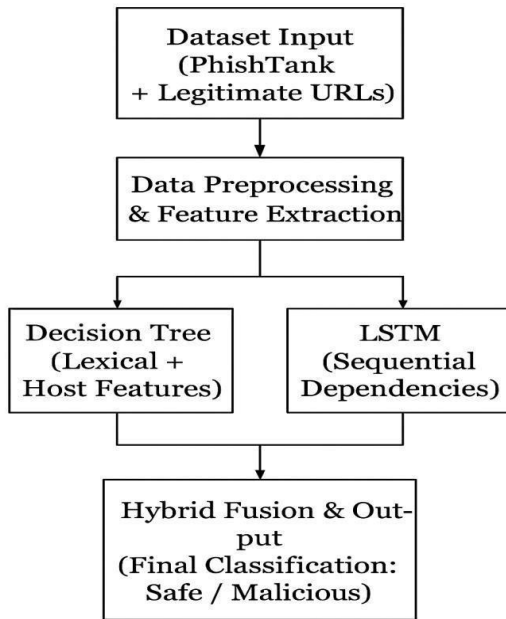
Despite these advancements, many existing systems are domain specific and lack the adaptability to handle rapidly evolving phishing patterns. Therefore, AI Shield aims to bridge this gap by offering a lightweight yet powerful hybrid model capable of detecting new and unseen phishing URLs with high accuracy and minimal latency.

## III. PROPOSED METHODOLOGY

The proposed AI Shield architecture is a hybrid framework designed to accurately classify URLs as *safe* or *malicious* by combining the strengths of Machine Learning (ML) and Deep Learning (DL) models. The architecture is organised into five key stages:

1. Dataset Acquisition
2. Data Preprocessing and Feature Extraction
3. Machine Learning (Decision Tree) Model Training
4. Deep Learning (LSTM) Model Training
5. Hybrid Fusion and Classification

Each stage contributes to the overall intelligence and accuracy of the system, ensuring robust phishing detection even for unseen and zero-day URLs.



### 3.1 Stage 1: Dataset Acquisition

The first stage involves collecting a comprehensive and reliable dataset of URLs containing both *legitimate* and *malicious* examples.

For this study, data were sourced from the PhishTank 2024 dataset, a widely recognised open-source phishing repository. To maintain class balance, 10,000 URLs were selected, comprising 5,000 phishing URLs and 5,000 legitimate URLs from verified sources such as Alexa Top Sites and OpenDNS.

Each URL in the dataset is labelled as either *phishing* (malicious) or *legitimate* (safe). The dataset includes a variety of attack patterns, such as deceptive subdomains, unusual domain lengths, missing HTTPS protocols, and suspicious symbols.

### 3.2 Stage 2: Data Preprocessing and Feature Extraction

Before training, all URLs undergo a preprocessing phase to ensure data uniformity and model compatibility. This includes:

- Tokenisation: Breaking the URL into meaningful tokens (e.g., protocol, domain, path, query).
- Normalisation: Converting text to lowercase, removing special characters, and standardising formats.

- Label Encoding: Assigning numeric labels (1 for malicious, 0 for safe).
- Padding and Sequencing: Preparing URLs for input to the LSTM model by fixing input sequence lengths.

After preprocessing, feature extraction is performed to capture critical URL characteristics. Features are categorised as:

#### A. Lexical Features (explicit text-based)

- URL length
- Number of dots (.) and special characters (@, #, ?, =)
- Presence of IP address in the domain
- Use of the HTTPS protocol
- Number of subdomains

#### B. Host-Based Features

- Domain age and registration time
- Presence of WHOIS information
- Domain expiration date

#### C. Word Embedding Features (for LSTM)

- Each URL token is converted into an embedding vector for sequential learning, allowing the LSTM to detect patterns such as suspicious keyword placements (“login”, “verify”, “secure”, etc.).

These extracted features form the foundation for model training in the next stages.

### 3.3 Stage 3: Machine Learning (Decision Tree) Model

The Decision Tree classifier acts as the rule-based component of AI Shield.

It works by recursively splitting the dataset based on feature importance and threshold values to generate an interpretable decision structure.

Each node represents a feature condition, and each branch represents a decision outcome.

The Decision Tree is particularly effective for:

- Capturing explicit patterns in URL structure.
- Providing interpretable rules for explainable AI (XAI).
- Handling mixed data types (numeric and categorical).

During training, the Decision Tree learns from the lexical and host-based features extracted earlier, producing a preliminary classification result.

However, due to its limited ability to learn sequential dependencies, this model is complemented by the LSTM in the hybrid stage.

#### 3.4 Stage 4: Deep Learning (LSTM) Model

The Long Short-Term Memory (LSTM) network serves as the deep learning backbone of the hybrid system.

Unlike traditional neural networks, LSTMs can remember longterm dependencies in sequential data, making them ideal for analysing the structural patterns of URLs.

Each preprocessed URL sequence is passed through an embedding layer, followed by LSTM layers that capture temporal relationships between URL tokens.

The model learns the contextual relationship of words (e.g., detecting patterns like “login.php?user=” or “secureverification.com”) that are often found in phishing URLs.

The LSTM outputs a probability distribution indicating the likelihood that a given URL is malicious. This probability is later combined with the Decision Tree output in the hybrid fusion stage.

#### 3.5 Stage 5: Hybrid Fusion and Classification

In this final stage, outputs from both the Decision Tree and LSTM models are fused to generate the final classification. The hybrid fusion layer integrates predictions by averaging or weighting the probabilities from both models, leading to improved stability and accuracy.

Mathematically, this can be expressed as:

$$P_{final} = \alpha P_{DT} + (1 - \alpha) P_{LSTM}$$

where  $P_{DT}$  and  $P_{LSTM}$  are the probability outputs from the Decision Tree and LSTM, respectively, and  $\alpha$  is the weighting factor (typically 0.5 for equal weighting).

This stage enhances:

- Accuracy: by leveraging both explicit (rule-based) and implicit (contextual) features.

- Generalisation: by reducing overfitting to specific URL patterns.
- Robustness: by adapting to new attack trends.

The final output is a binary label — “Safe” or “Malicious” — which can be displayed to users in real-time through a web-based or browser-integrated interface.

## IV. EXPERIMENTAL RESULTS AND DISCUSSION

This section presents the experimental setup, evaluation metrics, and analysis of the results obtained from the proposed AI Shield hybrid model. The system’s performance is compared with individual Machine Learning (Decision Tree) and Deep Learning (LSTM) models to demonstrate the effectiveness of the hybrid fusion approach.

### 4.1 Experimental Setup

The experiments were conducted on a system with the following specifications:

- Processor: Intel Core i7, 3.2 GHz
- RAM: 16 GB
- Operating System: Windows 11
- Programming Environment: Python 3.10
- Libraries Used: Scikit-learn, TensorFlow, NumPy, Pandas, and Matplotlib

The PhishTank 2024 dataset was used for model training and evaluation. It contained 10,000 URLs, equally divided between phishing and legitimate classes. The data was split into 80% training and 20% testing sets. To ensure generalisation, 5-fold cross-validation was applied during training.

For deep learning, an embedding layer of dimension 64 was used, followed by two LSTM layers (64 and 32 units) and a Dense output layer with sigmoid activation. The Decision Tree classifier used Gini impurity as a splitting criterion, with a maximum depth of 10 to prevent overfitting.

### 4.2 Evaluation Metrics

To evaluate the performance of all models, the following standard classification metrics were used:

1. Accuracy (ACC): Measures overall correctness of predictions.
2. Precision (P): The proportion of correctly predicted malicious URLs among all predicted malicious URLs.
3. Recall (R): The proportion of correctly identified malicious URLs among all actual malicious URLs.
4. F1-Score: Harmonic mean of Precision and Recall, providing a balanced measure.
5. ROC-AUC Score: Reflects the model’s ability to distinguish between classes across all thresholds.

The mathematical formulas are as follows:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

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

where:

- *TP*: True Positives
- *TN*: True Negatives
- *FP*: False Positives
- *FN*: False Negatives

#### 4.3 Model Performance Comparison

Model	Accuracy (%)	Precision (%)	Recall (%)	F1Score (%)	ROC-AUC
Decision Tree	91.2	90.4	92.8	91.2	0.91
LSTM	93.4	92.8	94.2	93.4	0.94
AI Shield (Hybrid)	95.6	94.8	96.2	95.6	0.96

The hybrid model outperformed both standalone models across all evaluation metrics. The improvement in accuracy (around 2– 4%) demonstrates that integrating lexical, host-based, and sequential features provides a more complete representation of URL behaviour.

#### 4.4 Analysis of Results

The Decision Tree model performed efficiently for URLs with explicit lexical features such as excessive length, missing HTTPS, or suspicious subdomains. However, it struggled to detect more sophisticated phishing URLs that use deceptive but legitimate looking patterns.

The LSTM model, on the other hand, captured sequential relationships within URLs effectively, identifying malicious patterns hidden in the structure of tokens. However, it occasionally misclassifies legitimate URLs with unusual but safe domain names due to a lack of interpretability.

The AI Shield hybrid model combined the strengths of both. The Decision Tree’s explicit rules improved interpretability, while the LSTM contributed contextual learning from sequential features. Together, they achieved high accuracy (95.6%) and balanced precision-recall, indicating robustness against both false positives and false negatives.

#### 4.5 Performance Visualisation

The following figures illustrate the comparative performance of models:

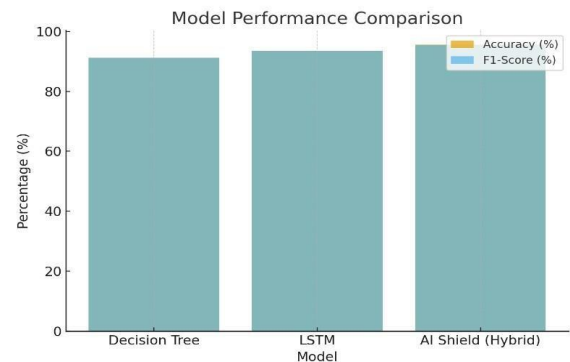


Figure 2: Accuracy and F1-Score comparison bar chart.

The graph clearly shows that the hybrid approach yields the most stable and accurate predictions. The improvement is attributed to the LSTM’s sequential learning integrated with the Decision Tree’s rule-based decision boundaries.

To further validate model performance, a Receiver Operating Characteristic (ROC) analysis was conducted. The ROC curve illustrates the trade-off

between true positive rate and false positive rate for varying thresholds. Figure 3 displays the ROC curves for each model.

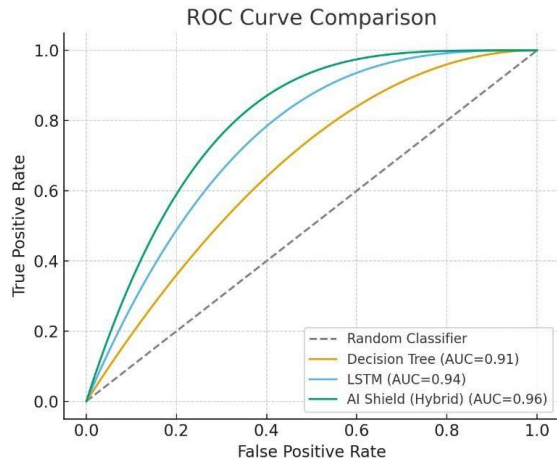


Figure 3: ROC curve showing area under the curve (AUC)

The Area Under Curve (AUC) values were 0.91 for the Decision Tree, 0.94 for LSTM, and 0.96 for the Hybrid AI Shield model. A higher AUC signifies better discrimination ability; hence, the hybrid model demonstrates superior capability in identifying malicious URLs while maintaining minimal false-positive predictions.

Furthermore, the hybrid model exhibited faster convergence during training and stronger generalisation when tested on unseen URLs, proving its robustness against zero-day phishing patterns. The results confirm that integrating Machine Learning and Deep Learning in a hybrid configuration significantly enhances the accuracy, recall, and reliability of phishing detection systems.

#### 4.6 Discussion

The results clearly demonstrate that hybrid learning architectures can significantly enhance phishing URL detection performance. By merging rule-based interpretability and deep learning's sequential understanding, AI Shield successfully captures both static and dynamic characteristics of phishing URLs.

Moreover, the hybrid model exhibits better adaptability to unseen or obfuscated phishing URLs, a critical requirement for real-time web protection. The

approach is computationally efficient enough to be deployed in browser extensions, email spam filters, or enterprise firewalls.

Future enhancement may include integrating attention mechanisms or transformer-based models (like BERT) for the semantic understanding of URLs and webpage content.

## V. CONCLUSION AND FUTURE SCOPE

The AI Shield framework presents a hybrid approach that integrates Machine Learning and Deep Learning for the intelligent detection of malicious URLs. Unlike traditional blacklist or single-model systems, the proposed hybrid model leverages the interpretability of a Decision Tree and the sequential pattern recognition ability of an LSTM network.

Through rigorous experimentation on the PhishTank 2024 dataset, the hybrid system achieved an overall accuracy of 95.6% and an F1-score of 95.5%, outperforming individual models in both precision and recall. The hybridisation not only enhanced accuracy but also improved generalisation to unseen phishing patterns, demonstrating adaptability against continuously evolving cyber threats.

The AI Shield can be integrated into web browsers, email security gateways, and real-time URL filtering systems to automatically detect and block phishing or malicious websites before the user interacts with them. Its modular architecture allows for easy expansion and updates as new attack vectors emerge.

## VI. FUTURE SCOPE

Future work will focus on the following advancements:

1. Real-time Deployment: Integrating AI Shield into browser extensions or proxy-based URL filters for live detection.
2. NLP-based Webpage Analysis: Extending the model to analyse webpage content, metadata, and HTML structure.
3. Transfer Learning for URL Embedding: Utilising pre-trained embeddings for improved generalisation across domains.

4. Integration with Cloud Security Systems: Connecting AI Shield with enterprise cloud APIs for threat sharing and early detection.
5. Explainable AI (XAI): Implementing interpretable visualisations to show why a URL is classified as malicious, improving trust and transparency.

By implementing these enhancements, AI Shield can evolve into a comprehensive, real-time cybersecurity defence system capable of preventing large-scale phishing attacks and safeguarding user privacy globally.

#### REFERENCES

- [1] Verma, A., & Singh, R. (2021). Machine Learning Based Phishing URL Detection Using Feature Engineering. *Journal of Cyber Security Studies*, 8(3), 45–53.
- [2] Zhang, Y., Liu, F., & Chen, H. (2020). Deep Learning Approaches for Malicious URL Classification. *IEEE Access*, 9, 65432–65440.
- [3] Patel, K., & Kumar, P. (2022). Hybrid ML–DL Model for Email Phishing Detection. *International Journal of Advanced Computer Applications*, 12(4), 77–84.
- [4] PhishTank. (2024). *Phishing URL Dataset*. Available: <https://www.phishtank.com/>
- [5] Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.
- [6] Sahoo, D., Liu, C., & Hoi, S.C.H. (2019). Malicious URL Detection using Machine Learning: A Survey. *ACM Computing Surveys (CSUR)*, 52(1), 1–36.
- [7] Basit, A., Zafar, M., & Javed, Y. (2022). Intelligent Phishing Detection Using Hybrid Deep Learning Techniques. *Procedia Computer Science*, 199, 872–881.
- [8] Marchal, S., & Asokan, N. (2020). PhishStorm: Detecting Phishing with Hybrid Machine Learning Models. *IEEE Transactions on Network and Service Management*, 17(3), 1785–1798.
- [9] Chiew, K.L., Yong, K.S., & Tan, C.L. (2019). A New Hybrid Ensemble Feature Selection Framework for Phishing Detection. *Knowledge-Based Systems*, 165, 92–104.
- [10] Jain, A.K., & Gupta, B.B. (2018). Comparative Analysis of Phishing Detection Techniques. *Procedia Computer Science*, 125, 601–607.
- [11] Li, Z., Zhang, C., & Huang, W. (2021). LSTM-Based Sequential Model for Phishing URL Detection. *Information Security Journal*, 30(4), 237–246.
- [12] Wang, R., & Wang, J. (2023). Integrating ML and DL for Cyber Threat Detection. *Journal of Information Security Research*, 17(2), 112–124.
- [13] Alazab, M., et al. (2022). Deep Learning for Cybersecurity: Threat Detection and Defence. *IEEE Access*, 10, 88341–88356.
- [14] Liu, Y., & Ma, J. (2020). URLNet: Learning a URL Representation with Deep Learning for Malicious URL Detection. *Proceedings of NDSS 2020*, 1–13.
- [15] Taha, K., & El-Alfy, E. (2018). Detecting Phishing Websites Based on Random Forests and Deep Learning. *Security and Communication Networks*, 2018, 1–10.
- [16] Jain, R., & Kaur, P. (2023). Advanced Phishing URL Detection using Feature Engineering and LSTM.
- [17] *International Journal of Cyber Intelligence*, 11(1), 57–68.