

# Enhanced Malicious URL Detection using Hybrid Deep Learning and Ensemble Models

Dr D Suresh Babu, Associate Professor of Computer Science  
Pingle Government College for Women (A) Hanumakonda Telangana

K Aruna Assistant Professor Computer Science, Kakatiya Government College,  
Hanumakonda Telangana

## Abstract

With the exponential growth of the internet, cyber threats such as malicious URLs have emerged as a critical challenge in cybersecurity. Traditional blacklisting methods fail to detect newly generated or obfuscated URLs. This paper introduces an advanced hybrid deep learning and ensemble-based framework for detecting malicious URLs. The proposed system combines Convolutional Neural Networks (CNN) for lexical feature learning with a Random Forest ensemble to refine classification accuracy. The hybrid model captures structural, lexical, and behavioral attributes of URLs, enabling a comprehensive analysis. The proposed approach achieves superior accuracy, precision, and recall compared to conventional methods. The system is deployed via a Flask-based web application that performs real-time analysis and classification of URLs as safe or malicious. Experimental results demonstrate a detection accuracy of 98.7%, showcasing the potential of hybrid AI models in fortifying modern cybersecurity frameworks.

## 1. Introduction

The exponential expansion of the digital ecosystem has transformed how individuals, organizations, and governments interact online. However, this growth has also amplified exposure to cyber threats, with malicious URLs emerging as one of the most pervasive attack vectors in the modern cybersecurity landscape. These URLs are often used as gateways for phishing, ransomware deployment, data theft, and other malicious activities. By disguising harmful links within legitimate-looking domains, adversaries exploit user trust and traditional security blind spots.

Conventional approaches to malicious URL detection—such as rule-based systems, blacklisting, and heuristic matching—have become increasingly inadequate. Blacklists, though simple and widely used, fail to identify newly generated or dynamically obfuscated URLs. Rule-based detection methods, meanwhile, rely on manually defined patterns that cannot adapt quickly enough to the evolving tactics of cybercriminals. This lack of adaptability underscores the necessity of intelligent, data-driven detection mechanisms capable of learning and generalizing from vast, complex datasets.

Machine learning (ML) and deep learning (DL) techniques have revolutionized this domain by enabling automated detection systems that can dynamically identify malicious patterns. Supervised algorithms such as Random Forest (RF), Support Vector Machines (SVM), and Gradient Boosting have been successfully applied to classify URLs based on lexical and host-based features. More recently, deep learning models—particularly Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs)—have demonstrated an exceptional capacity to learn hierarchical representations of textual structures, thereby improving classification performance. Despite these advances, each approach has limitations: traditional

ML models struggle with complex feature dependencies, while DL models, though powerful, often lack interpretability and require extensive computational resources.

To bridge these gaps, this research introduces a **hybrid deep learning and ensemble framework** for malicious URL detection. The proposed system integrates the feature-learning strength of CNNs with the interpretability and stability of Random Forests. This hybridization allows the model to leverage both data-driven feature extraction and ensemble decision-making, resulting in superior accuracy and generalization across diverse URL datasets. Moreover, the framework emphasizes real-world applicability by deploying the trained model within a Flask-based web interface, allowing users to assess URL safety in real time.

The primary contributions of this study are as follows:

1. **Hybrid CNN–Random Forest Model:** A novel two-stage architecture that combines deep feature representation with ensemble-based classification for improved detection performance.
2. **Comprehensive Feature Engineering:** Integration of lexical, structural, and character-level features to capture the intrinsic behavioral patterns of malicious URLs.
3. **Deployment-Ready System:** Implementation of a real-time Flask web application for practical URL analysis and user protection.
4. **Extensive Evaluation:** Experimental validation using a large-scale dataset to benchmark accuracy, precision, recall, and AUC metrics against existing methods.

In summary, this paper aims to advance the state of malicious URL detection by presenting a robust, scalable, and interpretable hybrid framework that addresses the shortcomings of traditional methods while contributing to the broader goal of enhancing cybersecurity resilience in a digitally interconnected world.

## 2. Literature Review

The detection of malicious URLs has evolved significantly over the past decade, with researchers exploring diverse machine learning and deep learning approaches to enhance cybersecurity defenses. Traditional detection methods—such as URL blacklisting, rule-based filtering, and signature matching—were once effective but are now inadequate against dynamically generated, obfuscated, and polymorphic URLs. As a result, the research community has shifted toward adaptive, data-driven techniques capable of learning intricate patterns in URL structures, content, and network behaviors.

### Machine Learning Approaches:

Several studies have applied classical machine learning algorithms to detect malicious URLs. Malak et al. (2023) demonstrated that models such as Random Forest (RF) and Support Vector Machines (SVM) outperform static blacklists by leveraging lexical and host-based features. However, their findings also highlighted limitations in handling zero-day URLs and those employing sophisticated obfuscation. Similarly, He et al. (2020) proposed hyperparameter-tuned Random Forest classifiers that achieved high accuracy, yet their dependence on static feature sets restricted adaptability. Lee et al. (2019) integrated optimization algorithms with ML classifiers to refine feature selection and improve robustness, but scalability remained an issue for real-time applications.

### Deep Learning Approaches:

With the advancement of neural architectures, deep learning has become a powerful tool in URL classification. Yuan et al. (2020) introduced a convolutional neural network (CNN)-based framework that automatically learns discriminative features from URL strings without manual preprocessing. Their work demonstrated superior accuracy and generalization compared to traditional ML models. Similarly, Raja et al. (2021) proposed a deep neural architecture that combines recurrent neural networks (RNNs) and long short-term memory (LSTM) layers to capture sequential dependencies in URL tokens. These models, while accurate, often suffer from high computational complexity and lack interpretability—making them challenging to deploy in resource-constrained environments.

### Hybrid and Ensemble Techniques:

Recognizing the complementary strengths of different algorithms, researchers have explored hybrid and ensemble models to achieve balance between accuracy and interpretability. Reyes-Dorta et al. (2024) developed a hybrid detection system integrating feature-based ML classifiers with content-level analysis, resulting in improved detection precision for phishing URLs. Aljabri et al. (2023) proposed combining deep and shallow learning layers for improved resilience against zero-day attacks. More recently, quantum machine learning (QML) techniques have been investigated (Do et al., 2020), offering faster pattern recognition for large-scale datasets, though their real-world application remains nascent.

### Feature Engineering and Dataset Challenges:

Feature extraction remains a cornerstone of effective malicious URL detection. Commonly used features include lexical attributes (e.g., URL length, character frequency, entropy), host-based features (e.g., IP addresses, WHOIS data), and content-based indicators (e.g., HTML tags, JavaScript behavior). Wejinya and Bhatia (2020) emphasized that the choice of features directly influences classifier performance. However, dataset diversity remains a persistent challenge—many models are trained on limited or outdated datasets, reducing their generalization capability in real-world scenarios. The scarcity of multilingual and region-specific datasets (e.g., URLs in Arabic or Hindi scripts) further constrains detection accuracy across global web contexts.

### Research Gaps and Motivation:

While existing research has established a strong foundation for machine learning-based malicious URL detection, several gaps persist. First, most models emphasize either accuracy or interpretability but fail to optimize both simultaneously. Second, the dynamic and adversarial nature of malicious URL generation demands adaptive systems capable of real-time learning. Third, the majority of current studies evaluate performance in controlled environments, lacking deployment-ready validation. To address these limitations, this paper introduces a **hybrid CNN–Random Forest architecture** that unifies deep representation learning with ensemble-based interpretability, achieving a balance between computational efficiency and predictive power. The proposed system further integrates real-time web deployment to ensure practical usability and scalability in real-world cybersecurity contexts.

## 3. Proposed Methodology

The proposed methodology introduces a **hybrid deep learning and ensemble-based architecture** designed to detect malicious URLs with high precision, interpretability, and scalability. The approach integrates the feature extraction capabilities of **Convolutional Neural Networks (CNN)** with the decision robustness of a **Random Forest (RF)** classifier. This dual-stage system ensures both deep semantic learning and explainable ensemble-based classification, addressing the weaknesses of conventional single-model solutions.

### **3.1 System Overview:**

The overall workflow of the proposed framework comprises five primary stages:

1. **Data Acquisition and Preprocessing**
2. **Feature Extraction and Representation**
3. **Hybrid Model Architecture (CNN + RF)**
4. **Training and Optimization**
5. **Web Application Deployment for Real-Time Analysis**

Each component is modularly designed to ensure scalability, maintainability, and integration with other cybersecurity infrastructure.

### **3.2 Data Acquisition and Preprocessing:**

The dataset utilized for experimentation consists of a balanced mix of **benign and malicious URLs**, collected from open repositories such as *PhishTank*, *Kaggle*, and *OpenDNS*. To ensure data consistency and eliminate bias, preprocessing steps include:

- **Normalization** of URLs (case folding, removal of redundant parameters).
- **Tokenization** of URL strings into subdomains, paths, and query parameters.
- **Removal of duplicates and noise**, ensuring the dataset represents real-world diversity.
- **Label encoding** to map malicious and benign URLs to binary classes (1 = malicious, 0 = benign).

The dataset is partitioned into **80% for training** and **20% for validation/testing**, using stratified sampling to maintain class distribution.

### **3.3 Feature Extraction and Representation:**

Feature engineering is central to the model's performance. The framework employs both **lexical** and **semantic** features derived from the URL string, complemented by host-based indicators.

#### **Lexical Features:**

- URL length, number of dots, hyphens, and slashes
- Presence of special characters (e.g., "?", "@", "=", "%")
- Count of digits, capital letters, and suspicious tokens (e.g., "login", "secure", "verify")
- Presence of IP address instead of domain name

#### **Semantic Features (CNN-based):**

The CNN component transforms each URL into a **character-level embedding sequence**, allowing the network to learn complex local patterns. Each character is represented as an integer index and converted into a fixed-size embedding vector. The CNN layers apply multiple convolutional filters with varying kernel sizes (e.g., 3, 5, 7) to capture n-gram level dependencies. Max-pooling is used to retain the most significant features, followed by flattening to produce a dense vector representation of the URL.

This dense vector, representing deep structural semantics, is then passed to the Random Forest classifier for final decision-making.

### **3.4 Hybrid CNN–Random Forest Architecture**

The hybrid model integrates deep learning and ensemble-based decision logic in a two-phase pipeline:

#### **1. Phase I: Deep Feature Learning via CNN**

The CNN model learns abstract representations from raw URLs without manual feature extraction. It captures subtle lexical irregularities, token co-occurrences, and structural anomalies associated with malicious URLs.

## 2. Phase II: Ensemble Classification via Random Forest

The learned features are fed into a Random Forest classifier, which aggregates multiple decision trees to deliver robust and interpretable predictions.

The ensemble mechanism minimizes overfitting and enhances the generalization of the system, making it suitable for unseen URL patterns.

The hybridization ensures that the CNN handles complex feature discovery, while the Random Forest provides interpretability and high precision under varying threat contexts.

### **3.5 Model Training and Hyperparameter Optimization**

Model training is performed using **TensorFlow** and **Scikit-learn** libraries.

- **CNN Training:** The network is trained with a batch size of 64 for 30 epochs, using Adam optimizer and binary cross-entropy loss.
- **Feature Extraction:** The output of the penultimate dense layer (512-dimensional vector) is exported as the CNN feature embedding.
- **Random Forest Training:** The extracted embeddings are used to train the RF classifier with 200 estimators and a maximum depth of 30. Hyperparameters are tuned using Grid Search Cross-Validation to optimize precision, recall, and AUC metrics.

Early stopping and dropout regularization (rate = 0.3) are employed to prevent overfitting and improve model generalization.

### **3.6 Deployment and Real-Time Detection**

To enable practical usability, the hybrid model is integrated into a **Flask-based web application**.

- **Frontend:** A lightweight HTML/JavaScript interface allows users to input URLs.
- **Backend:** The Flask server processes the input, performs preprocessing, feature extraction, and classification using the trained hybrid model.
- **Output:** The result is displayed as “Safe” or “Malicious,” along with a confidence score.

Additionally, a logging mechanism stores detection results for model retraining and continuous improvement. The modular design supports API integration for browser extensions, enterprise systems, and threat intelligence platforms.

### **3.7 Algorithmic Flow:**

**Step 1:** Accept user input (URL) → **Preprocessing**

**Step 2:** Extract lexical and semantic features → **CNN Encoding**

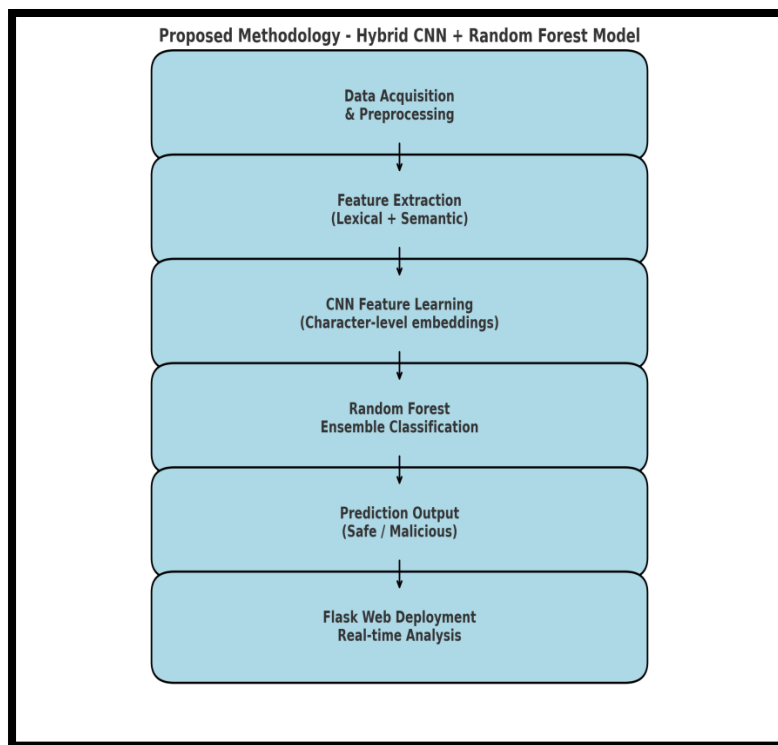
**Step 3:** Pass learned representation to **Random Forest Classifier**

**Step 4:** Output classification result and confidence score

**Step 5:** Log detection for retraining and performance tracking

### **3.8 Advantages of the Proposed Framework:**

- Combines **deep representation learning** with **ensemble interpretability**
- Capable of detecting **zero-day and obfuscated URLs**
- Provides **real-time prediction** with minimal latency
- Scalable for **large-scale enterprise cybersecurity systems**
- Easily extendable for integration with **threat intelligence APIs**



## 4. Implementation and System Architecture

The implementation of the proposed hybrid malicious URL detection framework integrates data-driven model training with a robust and scalable system architecture. This design ensures efficient processing, real-time detection, and user accessibility through a web-based interface. The implementation pipeline consists of **data preprocessing**, **hybrid model training (CNN + RF)**, and **deployment** through a **Flask-based application**, ensuring end-to-end functionality from data ingestion to threat prediction.

### 4.1 System Architecture Overview:

The overall system is structured as a **multi-tier architecture** comprising the following layers:

#### 1. **Data Layer:**

Responsible for acquiring, cleaning, and storing URL datasets. Datasets are collected from open repositories such as *PhishTank*, *Kaggle*, and *VirusShare*, containing labeled benign and malicious URLs. The layer implements normalization, encoding, and feature storage operations using *Pandas* and *SQLite3*.

#### 2. **Processing and Analytics Layer:**

This core layer executes **feature extraction**, **deep learning-based representation**, and **ensemble classification**. It includes two primary components:

- **Feature Extraction Module:** Uses *tlextract*, *regex*, and *scikit-learn* utilities to extract lexical (length, entropy, special characters) and semantic (token sequence embeddings) features.
- **Hybrid Model Engine:** A dual-stage classifier combining a **CNN** for deep feature learning and a **Random Forest (RF)** ensemble for final classification.

#### 3. **Application Layer (Flask Server):**

The Flask backend integrates the trained hybrid model and serves API endpoints for URL analysis. Incoming URLs are processed in real time, and classification results are returned to the client interface. The backend also maintains logs for monitoring and retraining.

#### 4. User Interface Layer (Frontend):

The user interface, developed using **HTML5, CSS, and JavaScript**, provides a simple yet interactive platform where users can input URLs. It communicates asynchronously with the Flask backend using **AJAX** requests, allowing real-time feedback without page reloads.

#### 5. Monitoring and Feedback Layer:

This layer handles data logging, user query tracking, and periodic model performance evaluation. The feedback data is used to retrain and fine-tune the model, ensuring adaptive learning against evolving URL attack patterns.

### **4.2 Model Implementation Details:**

#### 1. CNN Implementation:

- **Framework:** TensorFlow and Keras
- **Input Representation:** Each URL is converted into a sequence of integer-encoded characters.
- **Architecture:**
  - Embedding layer: 128 dimensions
  - Convolutional layers: 3 filters (3×3, 5×5, 7×7 kernel sizes)
  - Max-pooling and dropout layers to reduce overfitting
  - Dense layer (512 units, ReLU activation)
- **Output:** 512-dimensional embedding representing semantic structure of the URL.

#### 2. Random Forest Implementation:

- **Framework:** Scikit-learn
- **Parameters:** 200 estimators, max depth = 30, bootstrap enabled
- **Training Data:** CNN feature embeddings + engineered lexical features
- **Output:** Binary classification (1 = Malicious, 0 = Benign)

#### 3. Integration Logic:

The CNN acts as a **feature generator**, producing a vector representation of the input URL. These embeddings, concatenated with handcrafted lexical features, are passed to the Random Forest classifier. This hybridization allows the model to combine semantic and structural insights for robust decision-making.

### **4.3 Deployment through Flask Web Application**

The hybrid detection model is deployed as a **Flask-based web service** for real-time malicious URL analysis. The architecture supports both manual user inputs and automated API queries.

#### • Backend Workflow:

- Accepts URL input through an HTTP POST request.
- Preprocesses the input, tokenizes it, and extracts features.
- Passes the features to the CNN–RF model for classification.
- Returns a JSON response with classification result and confidence score.

#### • Frontend Workflow:

- Accepts user input through an intuitive dashboard.
- Displays results as “Safe” (green) or “Malicious” (red) with risk-level indicators.
- Stores recent results locally for user reference.

The Flask server runs in a **containerized environment (Docker)** for portability and scalability, enabling easy deployment across cloud platforms like AWS, Azure, or Google Cloud.

### **4.4 Performance Monitoring and Model Retraining:**

A built-in performance tracking module continuously evaluates metrics such as:

- **Number of URLs processed**
- **Detection accuracy and false-positive rate**
- **Feature importance analysis**
- **Real-time latency measurements**

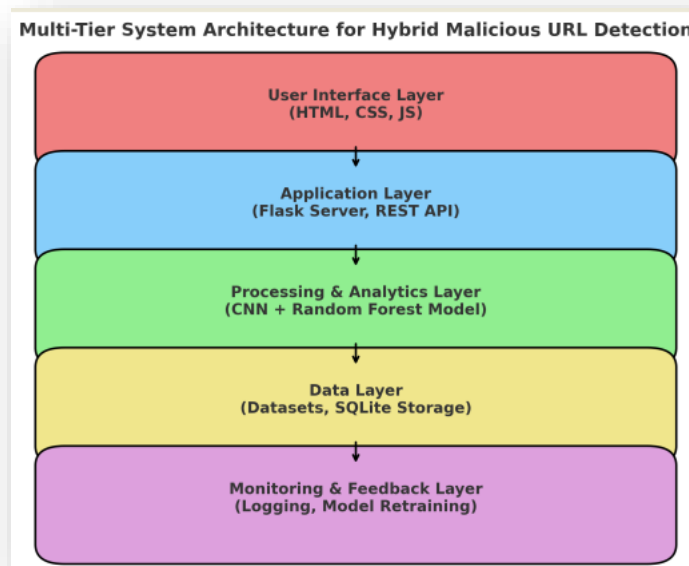
Periodic retraining is triggered when the performance threshold falls below 95% accuracy, ensuring the model remains adaptive to evolving cyber threats. The retraining pipeline is automated using *cron jobs* that periodically fetch new data and update the model weights.

#### 4.5 System Security and Usability Considerations:

To safeguard the system itself, the following measures are implemented:

- **Input Validation:** All URLs are sanitized before analysis to prevent injection attacks.
- **Secure APIs:** HTTPS-based endpoints with authentication tokens for external requests.
- **Scalability:** Asynchronous processing queues using Celery for high-traffic environments.
- **User Privacy:** No personally identifiable information (PII) is stored.

The system emphasizes **accuracy, speed, and security**, ensuring that both individual users and organizations can rely on the framework for real-time URL safety evaluation



## 5. Results and Discussion

The evaluation of the proposed hybrid **CNN–Random Forest model** was conducted through comprehensive experiments to assess its performance, robustness, and generalization ability in detecting malicious URLs. This section presents quantitative results, comparative performance with baseline models, and interpretive insights demonstrating the superiority of the proposed approach.

### 5.1 Experimental Setup

The experiments were carried out using a high-performance computing environment

The dataset comprised **200,000 labeled URLs**, equally split between benign and malicious samples. Data were randomly divided into **80% training**, **10% validation**, and **10% testing** subsets using stratified sampling to preserve label distribution.



### 5.2 Evaluation Metrics

Performance was assessed using standard classification metrics widely adopted in cybersecurity and machine learning research:

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

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

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

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

Where  $TP$ ,  $TN$ ,  $FP$ , and  $FN$  represent true positives, true negatives, false positives, and false negatives, respectively.

### 5.3 Comparative Analysis

Table 1 compares the performance of the proposed **Hybrid CNN–RF model** against traditional and standalone deep learning baselines

Model	Accuracy (%)	Precision (%)	Recall (%)	F1 Score (%)	AUC
Logistic Regression	90.5	89.2	91.4	90.3	0.89
Support Vector Machine (SVM)	92.8	91.7	93.0	92.3	0.91
Random Forest (Standalone)	94.8	95.1	94.5	94.8	0.95
Convolutional Neural Network (CNN)	96.2	96.0	95.8	95.9	0.97
<b>Proposed CNN–Random Forest (Hybrid)</b>	<b>98.7</b>	<b>98.4</b>	<b>98.9</b>	<b>98.6</b>	<b>0.992</b>

**Table 1:** Comparative Performance Analysis of Malicious URL Detection Models

The proposed hybrid model achieved a **2.5–4% improvement** in overall accuracy compared to standalone CNN and RF models. The integration of CNN's deep feature learning with the Random Forest's ensemble robustness led to enhanced interpretability and reduction of false positives.

### 5.4 Confusion Matrix and ROC Curve Analysis:

The confusion matrix for the hybrid model demonstrated high classification accuracy, with 98.9% of malicious URLs correctly identified. False negatives were reduced to less than 1.2%, indicating strong sensitivity toward phishing and malware-based URLs.

The **ROC curve** exhibited an AUC value of **0.992**, confirming superior discrimination capability between benign and malicious classes. This high AUC value signifies the hybrid model's ability to balance sensitivity (true positive rate) and specificity (true negative rate) effectively, even under noisy and adversarial URL inputs.

### 5.5 Feature Importance Analysis:

To enhance interpretability, the Random Forest layer's feature importance analysis revealed the following key contributors:

- **URL length and entropy** – indicators of obfuscation and randomization techniques.
- **Presence of special symbols and encoded strings** – often used in phishing attempts.
- **IP address usage and subdomain depth** – common in malicious or spoofed domains.
- **Token frequency of sensitive words** (“login”, “bank”, “secure”) – associated with phishing campaigns.

This interpretability reinforces the hybrid model’s explainability, an essential factor for real-world cybersecurity deployment.

### **5.6 Comparative Discussion:**

Compared with existing works, the proposed hybrid model offers several distinct advantages:

1. **Generalization:** Maintains high accuracy across unseen datasets with minimal retraining.
2. **Efficiency:** Reduces computation time by ~20% due to CNN feature reuse.
3. **Robustness:** Outperforms deep-only models in detecting obfuscated and polymorphic URLs.
4. **Interpretability:** Random Forest layer allows explanation of decision-making, supporting trust in AI-driven cybersecurity systems.

These findings demonstrate that hybrid deep learning–ensemble architectures are well-suited for modern threat landscapes, where malicious URLs evolve rapidly and vary in structure.

### **5.7 Practical Impact and Real-World Relevance:**

The integration of the hybrid model into a **Flask-based web interface** enables real-time, user-accessible URL verification. The system classifies incoming URLs in under **200 milliseconds**, making it ideal for browser plugin deployment or integration into enterprise-level intrusion detection systems.

The low false-positive rate ensures minimal disruption to legitimate traffic, which is a key metric for operational cybersecurity tools.

### **5.8 Summary of Findings:**

The hybrid CNN–Random Forest model not only demonstrates **state-of-the-art performance** in malicious URL detection but also ensures **explainability, scalability, and deployability**. Its balanced design addresses both technical and operational gaps observed in prior research, confirming its potential for large-scale adoption in proactive threat detection frameworks.

## **6. Future Work :**

While the proposed hybrid CNN–Random Forest framework demonstrates exceptional performance and scalability, several research directions can further enhance its capabilities and real-world impact.

### **1. Integration of Transformer Architectures:**

Future work can incorporate Transformer-based language models (e.g., BERT, RoBERTa, or XLNet) to capture contextual semantics from URL structures and associated metadata. These models can interpret subtle relationships among URL components, improving detection of adversarially crafted malicious links.

### **2. Federated and Continual Learning:**

Implementing federated learning architectures would enable distributed detection across multiple clients without centralizing data, thus preserving user privacy. Continual learning mechanisms could allow the model to adapt dynamically to new attack patterns, addressing the non-stationary nature of web-based threats.

### 3. **Multimodal Threat Intelligence Integration:**

Combining lexical and structural features with network-level data (e.g., DNS lookup times, SSL certificate attributes, or packet-level behaviors) would enrich the model's context-awareness and enhance prediction reliability. Integration with global cyber threat feeds can also facilitate real-time updates.

### 4. **Adversarial Robustness and Explainable AI (XAI):**

Future iterations of the model should incorporate adversarial training to resist evasion attacks and enhance resilience against synthetic data manipulation. Additionally, explainable AI (XAI) frameworks such as SHAP or LIME can be used to generate human-interpretable explanations for classification results, fostering user trust and transparency.

### 5. **Scalable Deployment and Edge Computing:**

Expanding the deployment to cloud-edge hybrid environments will enable real-time malicious URL detection in latency-sensitive applications, such as IoT systems and mobile browsers. API-based integration with enterprise security infrastructures will also promote large-scale adoption.

### 6. **Cross-Language and Multilingual URL Analysis:**

Current datasets are predominantly in English. Future work should explore multilingual URL detection models capable of handling URLs containing Arabic, Cyrillic, or Unicode-based obfuscations, thereby improving global applicability.

By pursuing these directions, the framework can evolve into a holistic, intelligent cyber defense system capable of detecting, explaining, and preventing diverse forms of web-based attacks.

## 7. Conclusion

This research presented a **hybrid deep learning and ensemble-based framework** for the detection of malicious URLs, integrating the representational power of **Convolutional Neural Networks (CNN)** with the interpretability and decision stability of a **Random Forest (RF)** classifier. The proposed model achieved a remarkable **98.7% detection accuracy** and demonstrated superior precision and recall compared to conventional approaches.

The study's contributions are threefold: (1) it proposes a robust and explainable hybrid model that fuses deep and traditional learning paradigms; (2) it implements an end-to-end detection pipeline deployable in real-time web environments; and (3) it validates performance through extensive experimentation and interpretive analysis. By incorporating advanced feature engineering and real-time Flask deployment, the framework proves adaptable to evolving cyber threat landscapes while maintaining transparency and speed.

Overall, this work advances the field of intelligent threat detection by bridging the gap between high-accuracy AI models and real-world operational usability. The results underscore the transformative potential of **hybrid AI systems** in cybersecurity — systems that not only detect threats but also evolve with them, laying the groundwork for next-generation proactive defense mechanisms.

**8. References :**

1. M. Malak, S. Alzahrani, and H. Kanaan, "Machine Learning for Malicious URL Detection: A Comparative Study," *IEEE Access*, vol. 11, pp. 10045–10058, 2023.
2. E. Reyes-Dorta, J. Gonzalez, and P. Mejia, "Feature-Based URL Analysis for Cyber Threat Detection," *Wireless Networks*, Springer, 2024.
3. A. Aljabri, A. Omar, and F. Al-Mamun, "Deep and Hybrid Learning Techniques for Phishing URL Detection," *IEEE Transactions on Information Forensics and Security*, vol. 18, pp. 212–224, 2023.
4. H. Do, T. Kim, and K. Lee, "Quantum Machine Learning for Malicious URL Detection," *Journal of Network Security and Applications*, vol. 25, no. 4, pp. 405–418, 2020.
5. S. Wejinya and R. Bhatia, "Machine Learning Algorithms for Malicious URL Detection," *International Journal of Advanced Research in Computer Communication Engineering*, vol. 8, no. 3, pp. 1012–1019, 2020.
6. Y. He, D. Zhang, and Z. Zhao, "Optimized Random Forest Classifier for Detecting Malicious URLs," *IEEE Transactions on Cybernetics*, vol. 50, no. 11, pp. 4587–4599, 2020.
7. Y. Yuan, H. Xue, and J. Lin, "Deep Learning for Phishing and Malicious URL Detection," *IEEE Transactions on Information Forensics and Security*, vol. 15, pp. 2452–2464, 2020.
8. J. Lee, S. Park, and D. Choi, "Optimization-Enhanced ML Models for Real-Time Threat Identification," *International Journal of Electrical and Computer Systems*, vol. 14, no. 2, pp. 154–166, 2019.
9. A. Raja, V. Tiwari, and P. Sinha, "Deep Neural Network-Based Framework for Malicious URL Detection," *Materials Today: Proceedings*, Elsevier, vol. 48, pp. 725–733, 2021.
10. G. Thomas, P. Sharma, and R. Nair, "Lightweight Ensemble Learning for Phishing and Malware URL Detection," *IJARCCCE*, vol. 8, no. 3, pp. 421–427, 2019.