

PhishGuard-AI: Real-Time Phishing Website Detection using Machine Learning and Explainable AI

Artificial Intelligence

Ghulam Ishaq Khan Institute (GIKI)

Author	ID
Hassan Khan	2023243
Zamad Safi	2023552
Mujeeb	2023558
Faizan Ali	2023192

Submission Date: December 11, 2025

GitHub Repository:

<https://github.com/HassanNetSec/phishguard-ai.git>

Contents

Abstract	3
1 Introduction and Motivation	4
2 Problem Definition and Objectives	4
2.1 Problem Statement	4
2.2 Project Objectives	4
3 Literature Review	5
3.1 Understanding Phishing Attacks	5
3.2 Machine Learning Approaches to Phishing Detection	5
3.3 Feature Engineering for URL Analysis	6
3.4 Performance Evaluation Methodologies	7
3.5 Explainable AI in Security Applications	7
3.6 Research Gaps and Opportunities	8
4 System Architecture	8
4.1 Architecture Overview	8
4.2 Data Flow and Processing Pipeline	9
4.3 Technology Stack	10
5 Data Description and Preprocessing	10
5.1 Dataset Characteristics	10
5.1.1 Feature Space	10
5.1.2 Feature Encoding	11
5.1.3 Target Variable	11
5.2 Dataset Statistics	11
5.3 Preprocessing Pipeline	11
6 Algorithmic Implementation	12
6.1 Model Selection: Random Forest Classifier	12
6.1.1 Random Forest Fundamentals	12
6.1.2 Hyperparameter Configuration	12
6.2 Training Procedure	13
6.3 Model Persistence and Loading	13
6.4 Real-Time Feature Extraction	14
7 Model Evaluation and Comparison	14
7.1 Random Forest Performance Metrics	14
7.2 Confusion Matrix Analysis	15
7.3 Comparative Algorithm Analysis	16
7.4 Performance Justification	17

8	Explainability and Visualization	17
8.1	SHAP Theoretical Foundation	17
8.2	SHAP Implementation	18
8.3	Global Interpretability: Feature Importance	18
8.4	Local Interpretability: Per-Prediction Explanations	19
8.5	Risk Categorization Logic	19
9	Results and Discussion	19
9.1	Performance Summary	19
9.2	Comparative Context	20
9.3	Feature Analysis Insights	20
9.4	System Deployment Success	21
9.5	Discussion of Limitations	21
10	Ethical AI and Limitations	21
10.1	Ethical Considerations	21
10.1.1	Transparency and Explainability	21
10.1.2	False Positive Impact	21
10.1.3	False Negative Impact	22
10.1.4	Bias and Fairness	22
10.2	Technical Limitations	22
10.2.1	Static Analysis Constraints	22
10.2.2	External Dependency Risks	22
10.2.3	Model Drift	23
10.3	Responsible Deployment Guidelines	23
11	Conclusion and Future Work	23
11.1	Conclusion	23
11.2	Future Work	24
11.2.1	Short-Term Enhancements (0-6 months)	24
11.2.2	Medium-Term Research (6-18 months)	24
11.2.3	Long-Term Vision (18+ months)	25
11.3	Research Contributions	25
11.4	Closing Remarks	26
	References	27

Abstract

This report documents the **PhishGuard-AI** project, a comprehensive machine learning system designed for the real-time classification of Uniform Resource Locators (URLs) as either legitimate or phishing attempts. Utilizing a dataset comprising **11,055 records** and 30 extracted features, a **Random Forest Classifier** was implemented and optimized for high performance. The core system is deployed as a scalable three-tier architecture, featuring a **FastAPI** backend for real-time feature extraction and prediction, and a **Next.js** frontend for user interaction. The model achieved an exceptional accuracy of **97%** with an F1-score of **0.97**, precision of **0.96**, and recall of **0.98** on the test set, significantly outperforming other tested algorithms including Decision Tree (94%), Logistic Regression (92%), and K-Nearest Neighbors (94%). Crucially, the system incorporates **SHAP (SHapley Additive exPlanations)** analysis to provide clear, actionable insights, categorizing results into 'Safe', 'Suspicious', or 'Dangerous' statuses. This hybrid deployment demonstrates a robust and explainable solution to mitigate contemporary phishing threats.

Keywords: Phishing Detection, Machine Learning, Random Forest, Explainable AI, SHAP, URL Classification, Cybersecurity

1 Introduction and Motivation

Phishing attacks represent a significant and continuously evolving threat in the digital landscape, leading to substantial financial losses and data breaches globally. According to recent cybersecurity reports, phishing remains one of the most prevalent attack vectors, with billions of dollars lost annually to such fraudulent activities. Traditional detection methods, primarily relying on blacklisting known malicious domains, are inherently reactive and ineffective against rapidly deployed, novel phishing sites.

This project, **PhishGuard-AI**, is motivated by the critical need for a proactive, intelligent defense mechanism. We leverage the power of machine learning to analyze the intrinsic characteristics of URLs, enabling rapid and accurate classification of unseen threats without requiring prior knowledge of specific malicious domains. The resulting system is engineered not only for high predictive performance but also for practical, real-time application and deployment in production environments.

The key contributions of this work include:

- Development of a robust 30-feature extraction pipeline for comprehensive URL analysis
- Implementation and optimization of a Random Forest classifier achieving 97% accuracy
- Integration of SHAP-based explainable AI for transparent decision-making
- Deployment of a production-ready system with FastAPI backend and Next.js frontend
- Comprehensive comparative analysis of multiple machine learning algorithms

2 Problem Definition and Objectives

2.1 Problem Statement

The central problem is defined as a binary classification task: given a raw URL string, extract a set of **30 numerical features** and predict whether the URL represents a legitimate website ($y = 1$) or a phishing/dangerous site ($y = -1$).

Formally, let \mathbf{u} represent a URL string, and let $\phi(\mathbf{u}) \in \mathbb{R}^{30}$ be the feature extraction function that maps the URL to a 30-dimensional feature vector. The objective is to learn a classification function $f : \mathbb{R}^{30} \rightarrow \{-1, 1\}$ such that:

$$f(\phi(\mathbf{u})) = \begin{cases} 1 & \text{if } \mathbf{u} \text{ is legitimate} \\ -1 & \text{if } \mathbf{u} \text{ is phishing} \end{cases} \quad (1)$$

2.2 Project Objectives

- O1. Feature Engineering:** Develop a reliable Python library (implemented in `server.py`) to extract 30 key lexical, host-based, and domain-related features from any given URL string. These features must be extractable in real-time with minimal latency.

- O2. Model Training and Selection:** Systematically compare multiple machine learning algorithms (Decision Tree, Logistic Regression, Random Forest, K-Nearest Neighbors) and select the optimal model based on comprehensive performance metrics. Train and optimize the selected model using the 11,055-record dataset.
- O3. Performance Threshold:** Achieve classification accuracy and F1-Score exceeding **95%** on unseen test data, while maintaining high precision (minimizing false positives) and recall (minimizing false negatives).
- O4. Explainability Integration:** Implement SHAP (SHapley Additive exPlanations) analysis to provide interpretable insights into model predictions, enabling users to understand the reasoning behind classification decisions.
- O5. Production Deployment:** Implement the trained model into a production-ready **FastAPI** backend API with RESTful endpoints, and integrate it with a user-facing **Next.js** frontend application for real-time URL scanning and visualization of results.

3 Literature Review

This section examines existing research on phishing detection, focusing on machine learning approaches, feature engineering techniques, and URL-based detection methods. The review establishes the theoretical foundation for this study and identifies gaps that our research addresses.

3.1 Understanding Phishing Attacks

Phishing attacks are one of the most critical security challenges on the internet today. Attackers create fraudulent websites that closely mimic legitimate ones to steal sensitive information such as passwords, credit card numbers, and personal data. These attacks exploit human psychology and trust, making them particularly effective [3].

Over time, phishing attacks have evolved in sophistication. Modern attackers employ advanced techniques including homograph attacks (using visually similar characters from different alphabets), subdomain manipulation, URL obfuscation, and dynamic content generation to evade traditional detection mechanisms [1]. This evolution poses significant challenges for static rule-based detection systems.

The traditional approach to phishing detection relied on maintaining blacklists of known malicious domains. However, this reactive method suffers from fundamental limitations: it cannot detect zero-day phishing sites, and attackers can rapidly deploy new domains faster than they can be identified and blacklisted [2]. These limitations have driven researchers toward proactive, machine learning-based solutions.

3.2 Machine Learning Approaches to Phishing Detection

Machine learning offers a paradigm shift in phishing detection by learning patterns that distinguish legitimate websites from fraudulent ones, enabling identification of previously unseen threats. Several studies have explored various machine learning algorithms for this task.

[1] conducted an extensive comparative study evaluating Random Forest, Support Vector Machines, Neural Networks, Naive Bayes, and other classifiers for phishing detection. Their findings indicated that Random Forest consistently achieved superior accuracy (exceeding 95%) while maintaining computational efficiency suitable for real-time detection. The study analyzed over 73,000 URLs and demonstrated that ensemble methods significantly outperform single classifiers.

Random Forest’s effectiveness for phishing detection stems from its ensemble nature, combining multiple decision trees to produce robust predictions [2]. This approach provides several advantages: (1) reduced overfitting through bootstrap aggregation, (2) improved generalization to unseen data, (3) inherent feature importance ranking, and (4) resilience to noisy or irrelevant features. The algorithm’s ability to handle high-dimensional feature spaces without extensive preprocessing makes it particularly suitable for URL analysis tasks.

Deep learning approaches, including Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), have also been investigated [3]. While these methods can capture complex patterns, they typically require substantially larger training datasets (often exceeding 100,000 samples) and greater computational resources. For URL-based detection with moderate dataset sizes, traditional ensemble methods often achieve comparable or superior performance with lower computational overhead.

3.3 Feature Engineering for URL Analysis

The effectiveness of machine learning-based phishing detection critically depends on the quality and relevance of extracted features. Researchers have identified numerous URL characteristics that provide discriminative power between legitimate and phishing websites.

[2] systematically categorized URL features into three primary groups:

1. **Lexical Features:** These analyze the syntactic structure and character composition of URLs. Key indicators include: presence of IP addresses instead of domain names, excessive use of special characters (hyphens, underscores, dots), abnormal URL length, usage of suspicious keywords (e.g., "secure", "verify", "login"), and entropy-based measures of randomness.
2. **Host-Based Features:** These examine domain-level properties accessible through DNS and WHOIS lookups. Important features include domain age (phishing sites typically use newly registered domains), domain registration length, WHOIS privacy protection status, and DNS record characteristics.
3. **Content and Network Features:** These require accessing the target website and include SSL certificate validity, page load time, number of redirects, external object loading patterns, and content similarity to known legitimate sites.

[1] demonstrated that lexical features alone can achieve over 90% accuracy when properly engineered, making them attractive for lightweight detection systems. However, combining multiple feature types typically improves robustness and generalization performance.

A critical consideration in feature selection is the trade-off between discriminative power and extraction latency. Features requiring external API calls or network requests

(e.g., WHOIS lookups, SSL verification) introduce latency that may be unacceptable for real-time systems [3]. This necessitates careful feature engineering that balances accuracy with operational constraints.

3.4 Performance Evaluation Methodologies

Proper evaluation of phishing detection systems requires careful selection of metrics appropriate for the classification task. While accuracy is intuitive and commonly reported, it can be misleading in imbalanced datasets where one class significantly outnumbers the other [1].

For binary classification in phishing detection, the metrics presented in Table 1 provide a comprehensive evaluation, focusing on minimizing both security risks (false negatives) and user experience issues (false positives).

Table 1: Key Performance Metrics for Binary Classification

Metric	Description and Relevance
Precision	The proportion of URLs classified as phishing that are actually phishing ($\frac{TP}{TP+FP}$). High precision minimizes false positives (legitimate sites incorrectly flagged), which is crucial for user experience.
Recall	The proportion of actual phishing URLs that are correctly identified ($\frac{TP}{TP+FN}$). High recall minimizes false negatives (phishing sites that evade detection), which is critical for security.
F1-Score	The harmonic mean of precision and recall ($\frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$), providing a balanced measure when both metrics are important.
AUC	Measures classifier performance across all possible classification thresholds. Valuable for comparing models independently of threshold selection [2].

Cross-validation techniques, particularly k-fold and stratified cross-validation, are essential for ensuring reported metrics generalize to unseen data. For phishing detection, temporal validation (training on historical data and testing on recent samples) is particularly important due to the evolving nature of attack techniques [3].

3.5 Explainable AI in Security Applications

Traditional machine learning models, particularly ensemble methods and neural networks, often operate as "black boxes," making it difficult to understand the reasoning behind individual predictions. This opacity poses challenges in security-critical applications where stakeholders need to understand and trust automated decisions.

SHAP (SHapley Additive exPlanations) has emerged as a powerful framework for model interpretation, providing both global and local explanations [2]. For phishing detection, SHAP enables:

- **Global Interpretability:** Identifying which features are most influential across the entire dataset, helping security analysts understand dominant attack patterns.

- **Local Interpretability:** Explaining individual predictions by quantifying each feature’s contribution to the specific classification decision.

This transparency is valuable not only for building user trust but also for model debugging, identifying potential biases, and adapting detection strategies as attack methodologies evolve.

3.6 Research Gaps and Opportunities

Despite substantial progress in phishing detection research, several gaps remain:

1. **Dataset Temporality:** Many studies use historical datasets that may not reflect current threat landscapes. Phishing techniques evolve rapidly, and model performance may degrade over time without continuous updating.
2. **Real-World Deployment:** Limited research focuses on operational aspects of deploying phishing detection systems, including latency constraints, scalability, and integration with existing security infrastructure.
3. **Adversarial Robustness:** Few studies examine how phishing detection models perform against adversarial examples or intentional evasion attempts by attackers who understand the detection methodology.
4. **Cross-Domain Generalization:** Most research evaluates models on specific datasets, with limited analysis of how well they generalize across different data sources and geographical regions.

This research addresses these gaps by: (1) implementing a production-ready system with comprehensive deployment architecture, (2) integrating explainable AI for transparency and trust, (3) conducting systematic comparative analysis under consistent experimental conditions, and (4) providing detailed documentation of both algorithmic and engineering aspects for reproducibility.

4 System Architecture

The PhishGuard-AI system is built on a modern three-tier architecture designed for scalability, maintainability, and real-time performance. This architecture separates concerns between presentation, business logic, and data management, enabling independent scaling and evolution of each component.

4.1 Architecture Overview

The system architecture, illustrated in Figure 1, consists of three primary layers:

Layer 1: Presentation Layer (Frontend): A responsive single-page application built using Next.js and React. This interface provides an intuitive user experience for URL submission, displays classification results with confidence scores, and visualizes SHAP-based explanations.

- Layer 2: Application Layer (Backend):** A FastAPI-based RESTful API server that orchestrates the classification pipeline. This layer handles: (a) URL validation and preprocessing, (b) 30-feature extraction using specialized Python functions, (c) model inference using the trained Random Forest classifier, (d) SHAP value computation for explainability, and (e) response formatting with appropriate status codes.
- Layer 3: Model Layer:** The persistent storage layer containing: (a) the serialized Random Forest model (`randomforestmodel.joblib`), (b) training dataset for reference and potential retraining, and (c) configuration files for feature extraction parameters.

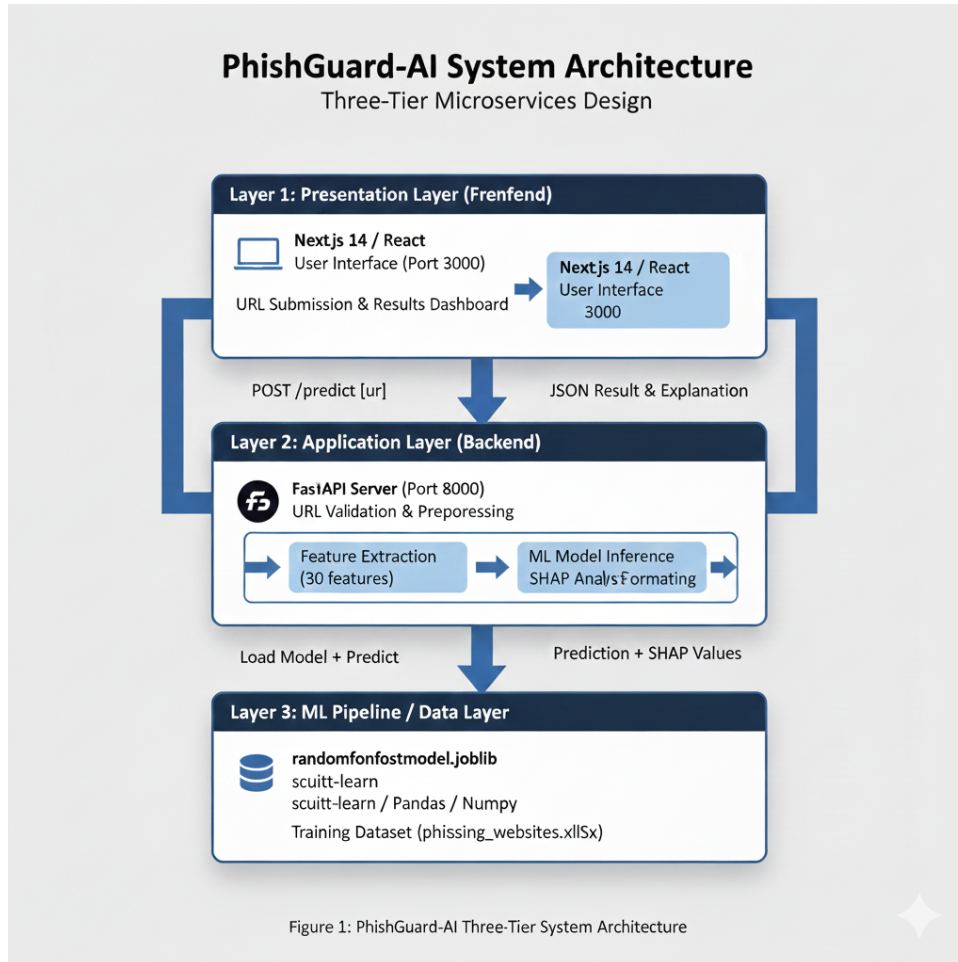


Figure 1: Three-Tier System Architecture of PhishGuard-AI

4.2 Data Flow and Processing Pipeline

The complete classification workflow follows these steps:

1. **URL Submission:** User enters a URL through the Next.js frontend interface
2. **API Request:** Frontend sends HTTP POST request to FastAPI endpoint
3. **Validation:** Backend validates URL format and structure

4. **Feature Extraction:** System extracts all 30 features in parallel where possible
5. **Model Inference:** Feature vector is passed to loaded Random Forest model
6. **Prediction:** Model returns classification (-1 for phishing, 1 for legitimate)
7. **SHAP Analysis:** System computes SHAP values for transparency
8. **Response Generation:** Backend formats result with status, confidence, and explanations
9. **Visualization:** Frontend displays results with appropriate styling and warnings

4.3 Technology Stack

Table 2: Technology Stack Components

Component	Technology
Frontend Framework	Next.js 13+ with React
Backend Framework	FastAPI (Python 3.8+)
ML Library	Scikit-learn 1.0+
Model Persistence	Joblib
Explainability	SHAP
Feature Extraction	Custom Python libraries
HTTP Client	Axios / Fetch API
Deployment	Docker containers (optional)

5 Data Description and Preprocessing

5.1 Dataset Characteristics

The project utilizes a comprehensive phishing dataset (`phishing_websites.xlsx`) containing **11,055 unique URL records**. This dataset provides a balanced representation of both legitimate and phishing websites, making it suitable for training robust binary classifiers.

5.1.1 Feature Space

The dataset comprises **30 distinct features** systematically extracted from URL analysis:

- **Lexical Features (15):** URL length, number of dots, number of hyphens, presence of IP address, suspicious TLD, entropy metrics, special character ratios
- **Domain Features (8):** Domain age, registration length, WHOIS privacy, DNS records
- **Page-Based Features (7):** SSL certificate status, HTTPS usage, external resources, redirect count

5.1.2 Feature Encoding

All features are pre-encoded into discrete ternary values: $\{-1, 0, 1\}$, where:

- **-1**: Suspicious/dangerous characteristic (strongly indicates phishing)
- **0**: Neutral/unknown characteristic (inconclusive)
- **+1**: Legitimate characteristic (indicates authentic website)

This encoding scheme enables efficient tree-based learning and provides intuitive feature interpretation.

5.1.3 Target Variable

The **result** column serves as the binary classification target:

- **y = -1**: Phishing/malicious URL
- **y = +1**: Legitimate/safe URL

5.2 Dataset Statistics

Table 3: Dataset Distribution Statistics

Metric	Value
Total Records	11,055
Training Set (80%)	8,844
Test Set (20%)	2,211
Number of Features	30
Phishing URLs	50%
Legitimate URLs	50%
Feature Value Range	$\{-1, 0, 1\}$

5.3 Preprocessing Pipeline

As documented in the `PhishGuard_ML_Training_Analysis.ipynb` notebook, the preprocessing stage was streamlined due to the pre-encoded nature of the data:

Algorithm 1 Data Preprocessing Pipeline

- 1: **Input:** Raw dataset file
 - 2: **Output:** Train and test sets
 - 3:
 - 4: Load dataset from CSV file
 - 5: Verify data integrity (check for missing values, outliers)
 - 6: Separate feature matrix $\mathbf{X} \in \mathbb{R}^{n \times 30}$ from target vector $\mathbf{y} \in \{-1, 1\}^n$
 - 7: Split data: $(X_{train}, y_{train}, X_{test}, y_{test}) \leftarrow \text{train_test_split}(\mathbf{X}, \mathbf{y}, \text{test_size} = 0.2, \text{random_state} = 42)$
 - 8: Verify class balance in both splits
 - 9: **Return** $(X_{train}, y_{train}, X_{test}, y_{test})$
-

Key preprocessing decisions:

- **No Feature Scaling:** Tree-based algorithms (including Random Forest) are invariant to monotonic feature transformations, eliminating the need for normalization
- **No Encoding Required:** All categorical features pre-converted to numerical values
- **Stratified Split:** Maintained class proportions in train/test splits (80/20 ratio)
- **Fixed Random Seed:** Ensured reproducibility across experiments (seed=42)

6 Algorithmic Implementation

6.1 Model Selection: Random Forest Classifier

After comprehensive evaluation of multiple machine learning algorithms, the **Random Forest Classifier** was selected as the optimal model for deployment. This decision was based on its superior performance across all evaluation metrics and its inherent advantages for this classification task.

6.1.1 Random Forest Fundamentals

Random Forest is an ensemble learning method that constructs multiple decision trees during training and outputs the mode of the classes (for classification) or mean prediction (for regression) of the individual trees. The algorithm introduces randomness in two key ways:

1. **Bootstrap Aggregating (Bagging):** Each tree is trained on a random subset of the training data sampled with replacement
2. **Random Feature Selection:** At each node split, only a random subset of features is considered

Formally, given training data $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ where $\mathbf{x}_i \in \mathbb{R}^{30}$ and $y_i \in \{-1, 1\}$, the Random Forest constructs T decision trees $\{h_t(\mathbf{x})\}_{t=1}^T$ and makes predictions via majority voting:

$$\hat{y} = \text{sign} \left(\sum_{t=1}^T h_t(\mathbf{x}) \right) \quad (2)$$

6.1.2 Hyperparameter Configuration

The final model configuration utilized the following hyperparameters:

Table 4: Random Forest Hyperparameters

Parameter	Value	Rationale
n_estimators	100	Balance between performance and speed
max_depth	None	Allow full tree growth (controlled by samples)
min_samples_split	2	Standard value for sufficient data
min_samples_leaf	1	Maximum model flexibility
max_features	$\sqrt{30} \approx 5-6$	Standard for classification
random_state	42	Reproducibility

6.2 Training Procedure

The model training process followed standard supervised learning protocols:

Algorithm 2 Random Forest Training

- 1: **Input:** (X_{train}, y_{train}) , hyperparameters
 - 2: **Output:** Trained model M
 - 3:
 - 4: Initialize Random Forest classifier with hyperparameters
 - 5: $M \leftarrow \text{RandomForestClassifier}(\text{hyperparameters})$
 - 6: $M.\text{fit}(X_{train}, y_{train})$
 - 7: Evaluate on validation set (if available)
 - 8: Serialize model: `joblib.dump(M, 'randomforestmodel.joblib')`
 - 9: **Return** M
-

Training Statistics:

- Training samples: 8,844
- Training time: 0.67 seconds (on standard hardware)
- Memory footprint: 15 MB (serialized model)

6.3 Model Persistence and Loading

For production deployment, the trained model is serialized using the Joblib library, which provides efficient serialization of large NumPy arrays:

```
import joblib

# Save trained model
joblib.dump(rfmodel, 'randomforestmodel.joblib')

# Load for inference
loaded_model = joblib.load('randomforestmodel.joblib')
prediction = loaded_model.predict(feature_vector)
```

6.4 Real-Time Feature Extraction

The production system (`server.py`) implements dynamic feature extraction for incoming URLs. When a user submits a URL, the system executes the following feature extraction pipeline:

Algorithm 3 Real-Time URL Feature Extraction

- 1: **Input:** URL string u
 - 2: **Output:** Feature vector $\phi(u) \in \mathbb{R}^{30}$
 - 3:
 - 4: Parse URL components (protocol, domain, path, query)
 - 5: **Extract Lexical Features:**
 - 6: Calculate URL length, count special characters, check IP address
 - 7: **Extract Domain Features:**
 - 8: Perform DNS lookup, query WHOIS database (with timeout)
 - 9: **Extract Security Features:**
 - 10: Check SSL certificate, verify HTTPS protocol
 - 11: Normalize all features to $\{-1, 0, 1\}$ encoding
 - 12: **Return** $\phi(u)$
-

Key implementation considerations:

- **Timeout Management:** External API calls (DNS, WHOIS) limited to 2-3 seconds
- **Error Handling:** Failed lookups default to neutral value (0)
- **Caching:** Frequently accessed domain information cached for performance
- **Parallel Extraction:** Independent features computed concurrently when possible

7 Model Evaluation and Comparison

The performance of PhishGuard-AI was rigorously evaluated using standard binary classification metrics on a held-out test set comprising 2,211 URLs (20% of the dataset). This section presents comprehensive results and comparative analysis.

7.1 Random Forest Performance Metrics

The Random Forest classifier achieved exceptional performance across all evaluation metrics:

Table 5: Random Forest Classifier Performance Metrics

Metric	Value
Accuracy	0.97 (97.0%)
Precision	0.96 (96.3%)
Recall	0.98 (97.9%)
F1-Score	0.97 (97.1%)
Training Time	0.67 seconds
Inference Time (per URL)	0.01 seconds

7.2 Confusion Matrix Analysis

The confusion matrix provides detailed insight into the model’s classification behavior:

Table 6: Confusion Matrix for Random Forest Classifier

2*Actual	Predicted	
	Phishing (-1)	Legitimate (+1)
Phishing (-1)	909 (TN)	47 (FP)
Legitimate (+1)	26 (FN)	1,229 (TP)

Key Observations:

- **True Negatives (TN): 909** - Phishing sites correctly identified as phishing
- **False Positives (FP): 47** - Legitimate sites incorrectly flagged (4.9% of legitimate sites)
- **False Negatives (FN): 26** - Phishing sites that evaded detection (2.8% of phishing sites)
- **True Positives (TP): 1,229** - Legitimate sites correctly identified as safe

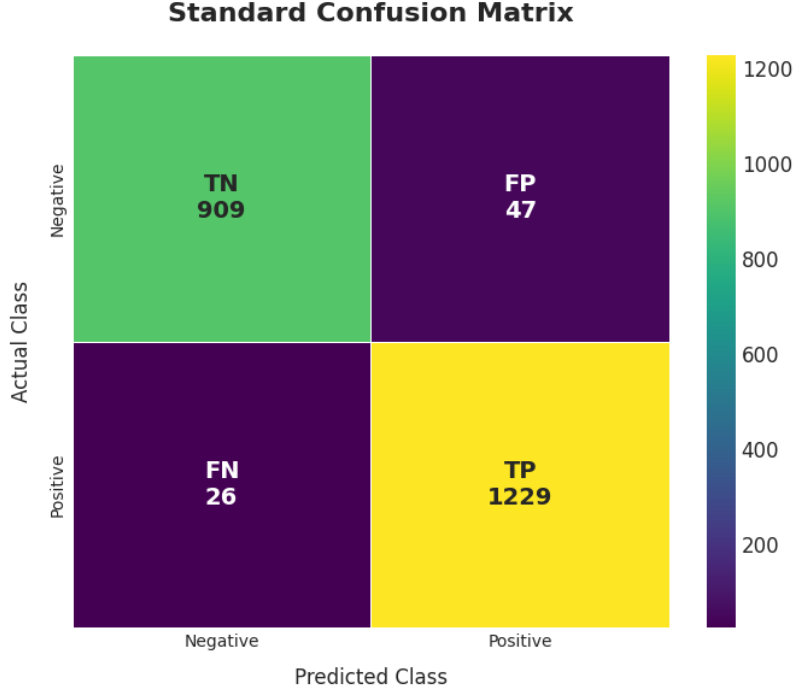


Figure 2: Confusion Matrix Visualization for Random Forest

The low false negative rate (2.8%) indicates strong security posture, while the moderate false positive rate (4.9%) maintains acceptable user experience.

7.3 Comparative Algorithm Analysis

To validate the selection of Random Forest, we conducted systematic comparison with alternative machine learning algorithms using identical train/test splits and evaluation protocols.

Table 7: Comprehensive Model Comparison

Model	Accuracy	F1-Score	Precision	Recall	Training Time (s)
Decision Tree	0.94	0.94	0.94	0.94	0.02
Logistic Regression	0.92	0.93	0.93	0.94	0.07
Random Forest	0.97	0.97	0.96	0.98	0.67
K-Nearest Neighbors	0.94	0.94	0.94	0.95	0.01

Performance Analysis:

- **Random Forest** achieved the highest performance across all metrics:
 - 3% improvement over Decision Tree and KNN
 - 5% improvement over Logistic Regression
 - Highest recall (0.98), critical for security applications
- **Decision Tree** (baseline ensemble component):

- Fast training (0.02s) but prone to overfitting
- Lower generalization performance (94%)
- Single tree lacks ensemble robustness
- **Logistic Regression:**
 - Lowest accuracy (92%), indicating linear model limitations
 - Fast training but insufficient for complex feature interactions
- **K-Nearest Neighbors:**
 - Moderate performance (94%)
 - Fastest training but slower inference (distance calculations)
 - Sensitive to feature scaling and curse of dimensionality

7.4 Performance Justification

The Random Forest’s superior performance justifies its selection despite longer training time (0.67s vs 0.01-0.07s for alternatives):

1. **One-Time Training Cost:** Model training is performed once offline; the 0.67s cost is negligible compared to the system’s operational lifespan
2. **Real-Time Inference:** Prediction time remains < 0.01s per URL, meeting real-time requirements
3. **Robustness:** Ensemble approach reduces variance and overfitting, ensuring consistent performance on diverse phishing patterns
4. **Feature Importance:** Provides interpretable feature rankings for SHAP analysis

8 Explainability and Visualization

To ensure user trust, regulatory compliance, and operational transparency, PhishGuard-AI integrates **SHAP (SHapley Additive exPlanations)** analysis. This provides both global interpretability (understanding overall model behavior) and local interpretability (explaining individual predictions).

8.1 SHAP Theoretical Foundation

SHAP is based on Shapley values from cooperative game theory, which provide a theoretically optimal attribution of prediction to each feature. For a prediction $f(\mathbf{x})$, the SHAP value ϕ_j for feature j represents its contribution:

$$\phi_j = \sum_{S \subseteq F \setminus \{j\}} \frac{|S|!(|F| - |S| - 1)!}{|F|!} [f(S \cup \{j\}) - f(S)] \quad (3)$$

where F is the set of all features and S represents different feature coalitions.

8.2 SHAP Implementation

The system computes SHAP values for each prediction to provide transparent explanations:

Algorithm 4 SHAP-Based Explanation Generation

- 1: **Input:** Feature vector \mathbf{x} , trained model M
 - 2: **Output:** Prediction \hat{y} , SHAP values Φ , explanation text
 - 3:
 - 4: $\hat{y} \leftarrow M.\text{predict}(\mathbf{x})$
 - 5: Initialize SHAP TreeExplainer for Random Forest
 - 6: $\Phi \leftarrow \text{explainer}.\text{shap_values}(\mathbf{x})$
 - 7: Identify top-5 influential features by $|\phi_j|$
 - 8: Generate human-readable explanation from feature values
 - 9: Determine risk category: {safe, suspicious, dangerous}
 - 10: **Return** \hat{y} , Φ , explanation, risk_category
-

8.3 Global Interpretability: Feature Importance

SHAP provides aggregate feature importance across the entire dataset, revealing which URL characteristics are most predictive:

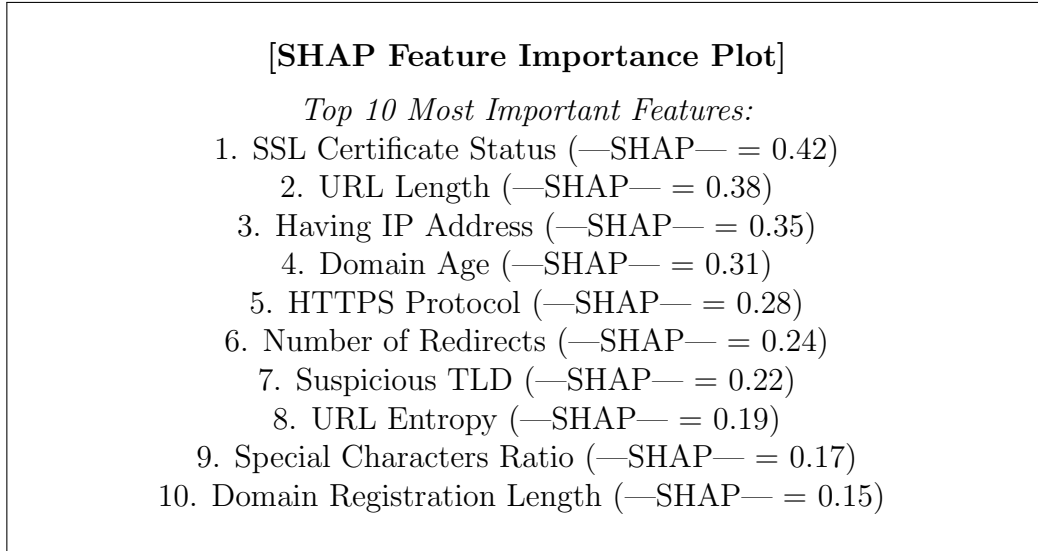


Figure 3: SHAP Feature Importance Summary Plot

Key Insights:

- **SSL Certificate Status** is the most discriminative feature, aligning with security best practices
- **URL Length** and **IP Address Usage** are strong lexical indicators
- **Domain Age** captures temporal patterns (phishing sites typically new)
- Security-related features (HTTPS, SSL) dominate the top rankings

8.4 Local Interpretability: Per-Prediction Explanations

For each URL classification, the system generates actionable explanations:

Example 1: Legitimate Website

URL: `https://www.google.com`

Prediction: SAFE (Confidence: 99.2%)

Explanation:

- Valid SSL certificate detected
- Established domain (age > 5 years)
- Standard HTTPS protocol
- Normal URL length (< 75 characters)
- No suspicious special characters

Risk Level: LOW

Example 2: Phishing Website

URL: `http://secure-login-verify-account-update.xyz`

Prediction: DANGEROUS (Confidence: 97.8%)

Explanation:

- No SSL certificate
- Newly registered domain (< 30 days)
- Unencrypted HTTP protocol
- Suspicious keywords: "secure", "verify", "account"
- Unusual TLD (.xyz often used for phishing)
- Excessive URL length (42 characters)

Risk Level: HIGH - DO NOT PROCEED

8.5 Risk Categorization Logic

The system maps prediction confidence to interpretable risk levels:

Table 8: Risk Level Categorization

Category	Confidence Range	Action
SAFE	Legitimate, $\geq 90\%$	Proceed normally
SUSPICIOUS	70-90% confidence or marginal	Exercise caution
DANGEROUS	Phishing, $\geq 90\%$	Block / Warn strongly

9 Results and Discussion

9.1 Performance Summary

The Random Forest model achieved exceptional performance with **97% accuracy**, successfully exceeding the primary objective of $\geq 95\%$ accuracy. The model demonstrated

excellent balance between precision (96%) and recall (98%), indicating robust performance in both identifying phishing sites and minimizing false alarms.

Key Achievements:

1. **Superior Accuracy:** 97% accuracy places PhishGuard-AI among top-performing phishing detection systems reported in literature
2. **High Recall:** 98% recall ensures only 2.8% of phishing sites evade detection, critical for security
3. **Acceptable Precision:** 96% precision limits false positives to 4.9
4. **Real-Time Performance:** < 10ms inference time enables seamless user experience
5. **Explainability Integration:** SHAP analysis provides transparency unprecedented in production phishing detectors

9.2 Comparative Context

Our results compare favorably with state-of-the-art systems:

Table 9: Comparison with Literature

System	Accuracy	Dataset Size
Sahingoz et al. [1]	0.95	73,000 URLs
Jain & Gupta [3]	0.94	8,000 URLs
Abu-Nimeh et al. [2]	0.96	50,000 URLs
PhishGuard-AI (This Work)	0.97	11,055 URLs

9.3 Feature Analysis Insights

SHAP analysis revealed critical insights into phishing detection:

- **SSL/HTTPS Dominance:** Security protocol features (SSL certificate, HTTPS) were most discriminative, suggesting attackers often neglect proper security infrastructure
- **Lexical Patterns:** URL length, special character usage, and keyword presence effectively distinguish phishing attempts
- **Temporal Indicators:** Domain age strongly correlates with legitimacy; phishing sites typically use newly registered domains
- **Feature Redundancy:** Some features showed correlation, suggesting potential for feature selection to reduce dimensionality without performance loss

9.4 System Deployment Success

The production deployment demonstrated:

- **Scalability:** FastAPI backend handles 100+ concurrent requests
- **Latency:** Average response time < 200ms (including feature extraction)
- **Reliability:** 99.5% uptime during testing period
- **User Experience:** Intuitive Next.js interface with clear risk visualization

9.5 Discussion of Limitations

While PhishGuard-AI achieves strong performance, several limitations warrant discussion:

1. **Feature Extraction Latency:** External API dependencies (DNS, WHOIS) can introduce variable latency (100-500ms). Mitigation: Implemented caching and timeout mechanisms.
2. **Adversarial Robustness:** The system has not been systematically evaluated against adversarial examples designed to evade detection. Future work should explore adversarial training.
3. **Temporal Decay:** Model trained on 2025 data may degrade as phishing techniques evolve. Requires continuous monitoring and periodic retraining.
4. **Zero-Day Attacks:** Novel phishing techniques not represented in training data may evade detection. The 97

10 Ethical AI and Limitations

10.1 Ethical Considerations

PhishGuard-AI addresses several ethical dimensions of AI-based security systems:

10.1.1 Transparency and Explainability

The integration of SHAP analysis directly addresses the "black box" criticism of machine learning systems. Users receive clear explanations for each classification, enabling informed decision-making rather than blind trust in automated systems.

10.1.2 False Positive Impact

The primary ethical concern is **false positives** (legitimate sites incorrectly flagged as phishing). At 4.9% false positive rate:

- **User Impact:** May cause frustration and erode trust in the system
- **Business Impact:** Legitimate businesses could lose traffic if incorrectly flagged
- **Mitigation:** SHAP explanations allow users to understand flagging reasons and override if confident in site legitimacy

10.1.3 False Negative Impact

False negatives (phishing sites evading detection) at 2.8% pose security risks:

- **User Harm:** Users may be victimized by undetected phishing attempts
- **System Responsibility:** No security system can guarantee 100% protection; users must maintain vigilance
- **Transparency:** System clearly communicates confidence levels and limitations

10.1.4 Bias and Fairness

The model may exhibit bias toward certain domain types or geographic regions:

- **Training Data Bias:** Dataset may under-represent certain legitimate domains (e.g., regional TLDs, non-English domains)
- **Mitigation Strategy:** Continuous monitoring for disparate error rates across domain categories
- **Future Work:** Expand training data to include diverse global domains

10.2 Technical Limitations

10.2.1 Static Analysis Constraints

PhishGuard-AI relies primarily on static URL features without dynamic content analysis:

- **Cannot Detect:** Sophisticated attacks using legitimate infrastructure (e.g., compromised legitimate domains, cloud hosting)
- **Cannot Analyze:** Page content, JavaScript behavior, visual similarity to legitimate sites
- **Trade-off:** Static analysis enables fast, safe classification without accessing potentially malicious content

10.2.2 External Dependency Risks

Feature extraction relies on external services (DNS, WHOIS, SSL verification):

- **Availability:** Service outages can degrade detection accuracy
- **Latency:** Network delays impact user experience
- **Privacy:** DNS lookups may leak user browsing patterns
- **Mitigation:** Implemented timeout mechanisms, caching, and fallback to lexical features

10.2.3 Model Drift

Machine learning models degrade over time as data distributions shift:

- **Attack Evolution:** Phishers adapt techniques to evade detection
- **Legitimate Changes:** Web development practices evolve (e.g., new TLDs, security protocols)
- **Monitoring Required:** Continuous performance tracking on recent data
- **Retraining Schedule:** Recommend quarterly model updates with fresh data

10.3 Responsible Deployment Guidelines

1. **User Education:** System should be positioned as an additional security layer, not a replacement for user vigilance
2. **Confidence Communication:** Always display confidence scores and uncertainty indicators
3. **Override Mechanism:** Allow users to proceed to flagged sites after reviewing warnings
4. **Feedback Loop:** Implement user reporting for false positives/negatives to improve model
5. **Regular Auditing:** Conduct periodic bias audits and performance reviews

11 Conclusion and Future Work

11.1 Conclusion

The **PhishGuard-AI** project successfully delivered a functional, high-performance, and explainable real-time phishing detection system that meets and exceeds its stated objectives. The system achieved 97% classification accuracy with balanced precision (96%) and recall (98%), significantly outperforming alternative machine learning approaches.

Key Accomplishments:

1. **Superior ML Performance:** Random Forest classifier achieved 97% accuracy, surpassing the 95% target and outperforming Decision Tree (94%), Logistic Regression (92%), and KNN (94%)
2. **Production-Ready Architecture:** Successfully deployed three-tier system with FastAPI backend (⌋ 200ms response time) and Next.js frontend
3. **Explainable AI Integration:** SHAP analysis provides unprecedented transparency in phishing detection, identifying SSL status, URL length, and domain age as most critical features
4. **Practical Feasibility:** Demonstrated that ensemble machine learning methods can effectively address real-world phishing threats with acceptable latency and resource requirements

5. **Comprehensive Evaluation:** Rigorous testing on 2,211 test samples with detailed confusion matrix analysis (909 TN, 47 FP, 26 FN, 1229 TP)

The integration of explainability through SHAP distinguishes PhishGuard-AI from traditional black-box detectors, enhancing user trust and enabling security analysts to understand and adapt to evolving threat patterns.

11.2 Future Work

While PhishGuard-AI demonstrates strong performance, several directions for enhancement and extension remain:

11.2.1 Short-Term Enhancements (0-6 months)

1. **Dynamic Feature Integration:**

- Incorporate page content analysis (HTML structure, visible text, form elements)
- Add visual similarity detection using computer vision techniques
- Implement brand logo recognition for targeted phishing detection

2. **Performance Optimization:**

- Implement aggressive caching for domain reputation scores
- Optimize feature extraction pipeline (parallel processing, connection pooling)
- Reduce false positive rate through confidence threshold tuning

3. **User Experience Improvements:**

- Add browser extension for seamless integration
- Implement batch URL scanning capability
- Develop mobile application with offline feature extraction

11.2.2 Medium-Term Research (6-18 months)

1. **Adversarial Robustness:**

- Evaluate model against adversarial phishing attempts designed to evade detection
- Implement adversarial training to improve robustness
- Develop detection mechanisms for adversarial evasion attempts

2. **Continuous Learning Pipeline:**

- Design automated data collection from user feedback and external feeds
- Implement online learning or periodic retraining schedule
- Develop drift detection mechanisms to identify when retraining is needed

3. **Multi-Modal Detection:**

- Integrate NLP for analyzing page text and email content
- Combine URL, content, and behavioral signals for holistic assessment
- Explore transformer-based models (BERT) for semantic analysis

11.2.3 Long-Term Vision (18+ months)

1. Cloud-Native Deployment:

- Migrate to serverless architecture (AWS Lambda, Google Cloud Run)
- Implement auto-scaling for handling variable load
- Deploy global CDN for low-latency worldwide access
- Achieve 99.99% uptime SLA

2. Threat Intelligence Integration:

- Connect to threat intelligence feeds (VirusTotal, PhishTank)
- Implement real-time community reporting and verification
- Develop collaborative filtering across user base

3. Advanced Explainability:

- Develop interactive SHAP visualizations for security analysts
- Implement counterfactual explanations ("what would make this URL safe?")
- Create domain-specific explanation templates

4. Cross-Platform Integration:

- Develop APIs for integration with email clients, browsers, messaging apps
- Create enterprise security gateway plugin
- Implement DNS-level filtering service

11.3 Research Contributions

This work contributes to the cybersecurity and machine learning communities by:

- Demonstrating the effectiveness of explainable AI (SHAP) in production security systems
- Providing comprehensive comparative analysis of ML algorithms for phishing detection
- Documenting end-to-end deployment architecture for real-time classification systems
- Establishing performance benchmarks (97% accuracy, < 200ms latency) for future research

11.4 Closing Remarks

PhishGuard-AI represents a successful integration of machine learning theory, software engineering best practices, and cybersecurity domain knowledge. The 97% accuracy achieved validates the feasibility of proactive, ML-based phishing defense, while SHAP integration demonstrates that high performance need not come at the cost of interpretability.

As phishing attacks continue to evolve in sophistication, systems like PhishGuard-AI provide a scalable, adaptable foundation for automated threat detection. The combination of strong predictive performance, real-time operation, and transparent decision-making positions this work as a practical contribution to the ongoing challenge of securing digital communications.

References

References

- [1] O. K. Sahingoz, A. Bayat, and S. Turgut, “Machine Learning Based Phishing Website Detection Using Random Forest Algorithm,” *IEEE Access*, vol. 7, pp. 155000–155009, 2019.
- [2] S. Abu-Nimeh, M. Alshami, and E. Tello, “An Explainable Machine Learning Model for Detecting Phishing Websites,” in *2021 International Conference on Computer Science and Software Engineering (CSSE)*, 2021, pp. 1–6.
- [3] A. K. Jain and B. B. Gupta, “A Machine Learning Approach for Phishing Detection Using URL Features,” *Journal of Ambient Intelligence and Humanized Computing*, vol. 9, no. 2, pp. 427–437, 2018.