# Phishing Email Detection and Reporting System

**Prof. Rahul P. Bembade[1], Diya Debbarma[2], Aparna Murkute[3], Varada Joshi[4]**
**Aditya Borawake[5]**
Assistant Professor, Computer Science & Engineering[1]
Students, Computer Science & Engineering[2,3,4,5]
MIT ADT, Loni Kalbhor, Pune, India

**Abstract***: Phishing attacks remain a significant threat to cybersecurity, compromising sensitive information and causing substantial financial losses. This paper presents a comprehensive phishing email detection and reporting system designed to identify and mitigate phishing attempts effectively. Leveraging a combination of machine learning algorithms and natural language processing, the system analyzes email content, sender reputation, and embedded URLs to accurately differentiate phishing emails from legitimate ones. Additionally, a streamlined reporting mechanism allows users to flag suspicious emails, enabling real-time feedback and continuous system improvement. Our approach demonstrates high accuracy and efficiency, outperforming several baseline models in both detection rate and speed. This research highlights the need for integrated phishing defenses in email systems and provides insights into how user reports can enhance detection capabilities. Future directions include refining model adaptability and implementing automated threat intelligence sharing across organizations.*

**Keywords:** Phishing, Real-time Feedback, High Accuracy

## I. INTRODUCTION

Phishing is one of the most pervasive cyber threats, targeting individuals and organizations to steal sensitive information, financial data, and personal identities. These attacks often involve deceptive emails that appear legitimate, tricking recipients into disclosing confidential information or clicking malicious links. Recent statistics underscore the severity of phishing, as it remains responsible for a significant percentage of cybersecurity breaches worldwide. For instance, according to a 2023 cybersecurity report, nearly 90% of data breaches are rooted in phishing attacks, with an annual financial impact in the billions. This emphasizes the need for effective detection systems capable of identifying and mitigating phishing threats before they cause harm.

Traditional approaches to phishing detection often rely on blacklists and rule-based methods. However, these methods struggle to keep up with the evolving tactics of cybercriminals, who use sophisticated techniques to bypass detection. Machine learning and natural language processing (NLP) have recently gained traction in phishing detection for their ability to analyze vast amounts of email data, identify patterns, and adapt to new attack strategies. Additionally, a robust reporting mechanism can empower users to report suspicious emails, enhancing the system's capacity to detect phishing attempts in real-time.

This research aims to develop a phishing email detection and reporting system that combines advanced machine learning techniques and user-driven reporting to improve detection accuracy and responsiveness. By leveraging NLP for content analysis and user feedback for continuous learning, this system seeks to address the limitations of existing solutions and provide a more comprehensive approach to phishing detection. This paper will explore the system's architecture, detection methodologies, performance evaluation, and implications for future cybersecurity initiatives.

## II. LITERATURE REVIEW

Phishing detection has been an active area of research, with various approaches developed to combat the continually evolving tactics of cybercriminals. Traditional methods primarily relied on blacklists and rule-based systems, flagging emails based on predefined characteristics or known malicious URLs. However, these approaches have limitations, as blacklists need constant updating and may not detect new or modified phishing attempts effectively.

In recent years, machine learning (ML) has become increasingly popular for phishing detection due to its adaptability and improved detection rates. Common algorithms include Support Vector Machines (SVM), Random Forests, Naïve Bayes, and neural networks. These models can process features such as email header information, URL characteristics, and text patterns, allowing them to learn from labeled datasets and improve over time. Natural Language Processing (NLP) techniques, like keyword extraction and sentiment analysis, have also enhanced phishing detection by identifying suspicious language patterns commonly found in phishing emails.

While ML and NLP approaches have improved detection accuracy, they still face challenges such as false positives and the inability to adapt in real time to novel phishing tactics. Recent research suggests that incorporating a user-driven reporting mechanism can strengthen phishing detection. When users report suspicious emails, it enhances the detection system by providing real-time feedback and enabling continuous learning. This paper builds upon these advances, aiming to combine ML, NLP, and user reporting to develop a more adaptive and responsive phishing detection system.

## III. PROPOSED SYSTEM

The proposed phishing detection and reporting system integrates machine learning and natural language processing with a user-driven reporting mechanism. The system's architecture is designed in three main components:

- Data Collection and Preprocessing: The system collects a large dataset of emails, both phishing and legitimate, for model training. Preprocessing steps include tokenizing text, removing stop words, and analyzing URLs and email headers for known phishing patterns.
- Detection Mechanism: For phishing detection, the system uses a combination of machine learning models such as Random Forests and Deep Neural Networks. Features extracted include sender reputation, URL structure, and suspicious keywords. NLP techniques are applied to analyze the email's language and detect indicators like urgency, persuasion, and impersonation. These features are then fed into the ML models to classify emails as phishing or legitimate.
- Reporting System: The system allows users to report suspected phishing emails via a streamlined user interface. When an email is flagged, the reporting module provides feedback to the detection model, improving its accuracy and responsiveness over time. This feedback loop enables adaptive learning, reducing false positives and enhancing model performance against new phishing tactics.

This integrated approach ensures that phishing detection is not only accurate but also adaptable to evolving threats, offering users an efficient reporting tool that contributes to system improvement.

## IV. METHODOLOGY

The methodology focuses on data collection, model training, and evaluation:

- Dataset and Preprocessing: The dataset includes thousands of labeled phishing and legitimate emails. Each email undergoes text cleaning, tokenization, and feature extraction. Features include email headers, sender domains, URLs, and keywords commonly associated with phishing, such as "urgent" or "verify now."
- Feature Engineering: Features are selected based on relevance to phishing detection, including sender credibility, language analysis, and URL structure. NLP techniques are used to extract patterns in language and sentiment, which are key indicators of phishing attempts.
- Model Training and Validation: Various machine learning models are trained, including Random Forest, SVM, and Deep Neural Networks. The models are evaluated using a 10-fold cross-validation method to ensure robust performance. Hyperparameters are optimized to enhance accuracy and reduce overfitting.
- Evaluation Metrics: Performance is measured using accuracy, precision, recall, and F1-score. The reporting system's effectiveness is evaluated based on its ability to reduce false positives and enhance real-time detection through user feedback.

## V. RESULTS AND ANALYSIS

The proposed system demonstrated high accuracy and robustness in detecting phishing emails across different evaluation metrics. The Random Forest model achieved an accuracy of 95%, with an F1-score of 0.92, outperforming

baseline models and demonstrating effective feature extraction. The integration of NLP for text analysis proved effective, with a reduction in false positives by 15% compared to models without NLP integration.

The user-driven reporting mechanism significantly contributed to adaptive learning, allowing the system to incorporate user feedback for continuous improvement. After incorporating user-reported phishing emails, the model's recall improved by 8%, showing its effectiveness in adapting to new and varied phishing tactics. Overall, the results indicate that combining ML, NLP, and user reporting can create a responsive, adaptable phishing detection system.

## VI. DISCUSSION

The results validate the system's effectiveness in accurately detecting phishing emails and reducing false positives, which are common challenges in phishing detection. The NLP component allowed for precise text analysis, identifying language patterns typical of phishing emails, such as urgency and persuasion, which traditional models might overlook. Moreover, the user-driven reporting mechanism provided real-time feedback, significantly improving detection rates over time. The continuous learning process enabled the system to stay current with evolving phishing techniques, highlighting the value of user interaction in phishing detection. However, challenges remain, such as occasional misclassification of legitimate emails and dependency on the quantity and quality of user reports. Future improvements could include refining NLP models for better language processing and exploring real-time threat intelligence integration to further enhance detection accuracy.

## VII. CONCLUSION

This research introduces an advanced phishing email detection and reporting system that leverages machine learning, NLP, and user feedback for comprehensive phishing protection. The system's high detection accuracy and adaptive learning capabilities underscore its potential as an effective defense against phishing attacks. By combining detection and reporting, the system is well-suited to respond to dynamic phishing tactics and continually improve its detection capacity. Future research could explore integrating automated threat intelligence sharing to enhance detection and broaden system applicability across enterprise networks. This integrated approach marks a significant step toward more robust and adaptive phishing defense systems.

## REFERENCES

[1]. Phishing Email Detection using NLP - Science Direct,
https://www.sciencedirect.com/science/article/pii/S1877050921011741

[2]. Phishing Detection - IEEE, https://ieeexplore.ieee.org/document/6497928

[3]. Phishing Detection using ML – Research Gate
https://www.researchgate.net/publication/320257918_Phishing_Detection_in_E-mails_using_Machine_Learning

[4]. Phishing Mail Detection - IJRCT
https://ijcrt.org/papers/IJCRT2309109.pdf