# ShieldPhish: A Multi-Layered Phishing Detection System

Vaishnavi Chirawande[1], Aditi Gade[2], Shweta Ahire[3], Sharvari Jadhav[4]

[1]*Department of Computer Science and Engineering, Vishwakarma Institute of Technology, Pune, Maharashtra, India*

[2,3,4] *Department of CSE AIML, Vishwakarma Institute of Technology, Pune, Maharashtra, India*

*Abstract*—**Phishing is a continuing and advancing cyber threat that targets users across several means of communication. This report presents the design and implementation of ShieldPhish, a complete, multi-layered system designed to detect phishing attempts in invented URLs, emails, and SMS communication. ShieldPhish employs various feature extraction methods, including lexical, semantic (with transformer-based embeddings), and contextual analysis, in conjunction with strong machine learning models: a stacked ensemble of XGBoost and Logistic Regression for URLs and SMS, and a custom Multi-Input Neural Network for emails. ShieldPhish integrates real-time threat intelligence feeds as an initial defense, and it employs SHAP and LIME to provide explainable predictions to help users trust and understand the system. We discuss approaches and considerations for data acquisition, pre-processing, the architecture of the detection models, threat intelligence integration, explainability approaches, and finally, the deployment of the system across several interfaces: a Python package with CLIs, a Streamlit web app, and a browser extension for Chrome with proactive interception of navigation. Challenges encountered during development, such as heterogeneous data, complexity of feature engineering, instability of models, and complexity of environments/packaging, are discussed alongside proposed solutions. Performance evaluation shows high performance across the data modalities to suggest ShieldPhish as a protective system against contemporary phishing threats.**

*Index Terms*—**Phishing Detection, Cybersecurity, Machine Learning, Ensemble Learning, Neural Networks, Transformer Embeddings, Explainable AI (XAI), SHAP, LIME, Threat Intelligence, Feature Extraction, URL Analysis, Email Security, SMS Filtering, Browser Extension, Streamlit Web App**

## I. INTRODUCTION

Phishing attacks are on the rise in the digital ecosystem, which are campaigns that deceive users to provide sensitive information or take damaging action. They capitalize on human vulnerabilities and can be delivered through many means, especially through a malicious URL, email, and SMS (Smishing). While browser blocklists and email spam filters provide some protection, they are still vulnerable to new phishing tactics that use advanced obfuscation.

ShieldPhish is addressing this issue by creating a single detection system that can analyze the content of the three primary vectors described. We intend to build a powerful, adaptable, and transparent phishing detection instrument that is compliant with security best practices, and leverages current machine learning capabilities. It also provides an explanation, along with a prediction, for its decision.

What sets ShieldPhish apart is its single system architecture that supports multi-modal detection capability,

direct binding of different feature types (lexical, semantic, contextual) into specially designed model architectures (stacked ensemble, multi-input neural network), use of live threat intelligence as a first line of defense, and the giving of explainable AI insights through several deployment interlocks. It also encompasses a full-scale defense against the complex multi-faceted nature of phishing.

## II. LITERATURE REVIEW

In the last couple of years, research has been made on the usage of machine learning (ML), deep learning (DL), and transformer-based models for the detection of phishing in communication channels such as URLs, emails, and SMS

Omari (2023) experimentally compared seven ML algorithms (e.g., Logistic regression, Random forest, Gradient boosting) and demonstrated that ensemble methods outperform traditional classifiers for phishing

website detection, achieving high accuracy and robustness [1]. Similarly, Alzboon et al. (2024) emphasized that a neural network and a Random forest were more effective, resulting in an AUC of 0.994 and accuracy of 96.4%[2]. Besides, they found the selection of features to be quite significant, for example, URL length and domain age. These discoveries have an impact on the statement of Sharma et al. (2023), who mixed machine learning with natural language processing (NLP) to classify URLs and claimed greater than 95% accuracy, highlighting the relevance of semantic and lexical features of URLs for detecting new phishing attempts [2].

Phishing detection has been improved dramatically by deep learning techniques, especially transformer-based methods. For instance, the paper by Mahmud Dipto et al. (2023) examined the concept of explainable AI (XAI) with vision transformers to demonstrate improvement in explainability and improved performance in security-related classification challenges. The findings presented in this paper demonstrated that transformer-based models, such as BERT and its variants, are shown to be much better than conventional deep learning models to detect phishing emails, as they capture contextual semantics and even subtleties present in phishing content. The future direction of phishing detection systems beyond email would be to develop multi-modal phishing detection systems that combine text, visual, and metadata features, which we expect, would provide for a more robust and reliable detection system capable of making accurate predictions, via an email, SMS, or web-based platform.

Combining Threat Intelligence (TI) into detection systems has been a prominent trend because of the ability of these systems to adjust to the new strategies being adopted in phishing. The utilization of real-time TI feeds and machine learning (ML) and deep learning (DL) models enhances the ability to find zero-day phishing and resist adversarial tactics. Also, there have been a couple of recent studies that discuss the subject of explainability in phishing detection, where XAI techniques provide transparency and offer assistance in establishing trust with the end-users and analysts. It is definitely true for companies that utilize automated phishing detection in extremely risky situations where understanding is necessary for incident response and compliance.

All papers from 2023 to 2025 point to a definitive transition to hybrid and multi-modal approaches that use the advantages of ML, DL, transformers, TI, and XAI to deal with the ever-changing phishing attack threat. The route to the next generation of phishing detection tools is being opened by the deployment of advanced feature engineering, real-time intelligence, and interpretable models to facilitate rapid, scalable, and interpretable solutions can be deployed across various communication channels.

## III. RELATED WORK

For a long time, phishing detection has remained a significant challenge that various researchers have attempted to solve. In the early days, approaches were taken using heuristic rules that relied on the identification of suspicious words, URL patterns, and email headers. Then, Statistical Machine Learning models like Support Vector Machines (SVM), Naive Bayes, and Random Forests demonstrated improved results by learning patterns from the annotated data.

One of the noticeable changes owing to deep learning is the employment of Neural Networks especially on text data (email body, URLs paths) where RNN and CNN are often used. In fact, the latest generation of language models, including BERT and its successors, have been acknowledged as the best performing models at understanding the semantic meaning of text. They have also been shown to be efficient in identifying forgery in phishing emails and \ SMS.

Ensemble methods, which fuse multiple models, are lauded for their effectiveness in increasing robustness and accuracy. Specifically, stacked ensembles are exceptionally successful at combining different types of models or features, and involve a meta-model learning from predictions (or features) produced by base models.

While the bulk of research publications focuses on single information channels (for example, URL detection or email spam detection), few line the literature about multi-modal detection and practical explainability-based systems, hence the contribution of ShieldPhish is greater.

## IV. DATA ACQUISITION AND PREPROCESSING

A key step was acquiring datasets that included both phishing and non-phishing content for the modalities of the study. This involved reading data of different formats, cleaning the data, converting the labels to a binary phishing / benign classification, merging datasets, removing duplicates, and addressing class imbalance .

SOURCES:

- URLs: The datasets were retrieved by manually downloading and uploading from local files, and the corresponding parsers were written to account for their unique structures (e.g., feed comments and distinct columns from the CSV). The datasets include URLhaus daily dumps, OpenPhish feed, PhishStats score feed (which was parsed from CSV), Mendeley Phishing Dataset, and Alexa Top 1M (sampled for benign).
- Emails: Datasets consisted of different CSV-based email datasets (CEAS_08, Enron, Ling, Nazario, Nigerian_Fraud, SpamAssian). The main challenges in parsing were inconsistency of column headers (subject, body, label), different label values (0/1, spam/ham), and no explicit header columns (headers were mostly in the body or in the flattened CSV format missing). Custom pandas loading functions (load_emails_from_csv_with_header) were created, with explicit column mapping and label normalization.
- SMS: We used the UCI SMS Spam Collection (CSV) dataset along with other Kaggle sources (such as *Dataset_sms.csv*). These datasets were relatively straightforward but required some preprocessing — like renaming columns for consistency and mapping labels (for example, converting *"Smishing"* entries into the broader *"Phishing"* category).

## V. SYSTEM ARCHITECTURE

ShieldPhish is organized around a single Python package that contains the models, feature pipelines, and explainability utilities. Trained artifacts are stored separately and loaded at runtime. To make the system usable in different contexts, we built three interfaces that all call the same prediction functions, so there is one canonical inference path.
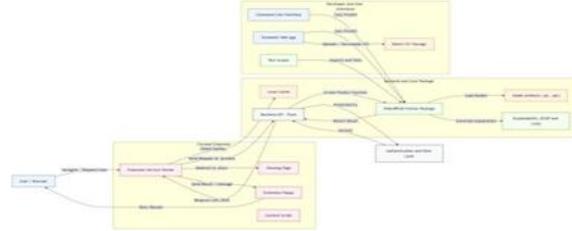


Figure 1. System Architecture Diagram

- Python package and command line interface: The core implementation lives in the shieldphish directory and is packaged for standard Python installation. The setup script allows both normal and editable installs with pip. A command line interface built with the click library exposes simple commands such as scanning for single inputs and checking for batch files. The CLI calls the same prediction functions used by the other interfaces, so behavior is consistent across tools.

- Streamlit web application: We implemented a Streamlit based web application for quick interactive use and demonstrations. The app provides a home page, a single scan page for paste input and instant explanations, a batch scan page that accepts a CSV and returns a results file, and a dashboard area reserved for future visualizations. The single scan uses a small heuristic to guess the input type and to extract subject and body when email text is provided. The batch scan expects a text column and a type column and can parse JSON encoded email content in CSV cells. The app reports parsing errors to the user rather than dropping rows silently.

- Backend API:
A minimal Flask API exposes a predicted endpoint and a placeholder report endpoint. The API accepts JSON payloads with the input type and text and returns prediction, score, and a compact explanation. The API reuses the package prediction functions, so the same code path is used whether clients call the CLI, the web app, or the browser extension. Production concerns such as authentication, rate limiting, TLS, and running behind a WSGI server are noted for future work

- Chrome extension:

The browser extension follows the Manifest V3 model and provides proactive navigation protection and on demand scanning. A service worker listens to navigation events and checks the destination URL using a local cache and the backend API. If the URL is classified as phishing above a conFigureurable threshold, the user is redirected to a custom warning page.

The warning page shows the scan result and offers concise actions: go back or continue to the site. Choosing continue temporarily bypasses checks for that URL. The extension also adds context menu items to scan selected text, links, or the current page and displays results in the popup. Manifest V3 required message passing and storage based communication between the service worker and the UI pages, and we tuned caching and thresholds to reduce repeated checks and avoid user frustration.

- Models, artifacts and explainability: Saved model artifacts include classic model pickles and neural network state files along with scalers and feature lists. At runtime the prediction module applies model specific preprocessing before inference. For explainability, we use SHAP for the tree-based models and a wrapper approach for the email neural network so interpretable attributions can be returned. All interfaces display a short human readable explanation that highlights the most influential features or words.

- Errors, limitations and recommended best practices:

Practical issues encountered during development include environment specific packaging problems and model artifact incompatibilities. The browser extension required careful handling of message passing because direct view access is not available in the Manifest V3 model. The URL model showed some aggressive false positives, so we used caching and a conFigureurable threshold to reduce interruptions for users. For production, we recommend hosting the API behind a WSGI server with TLS and authentication, adding rate limiting, and keeping model artifacts versioned and tied to the code that loads them.

ShieldPhish combines a single, canonical Python package with three user interfaces and a small backend so the same prediction logic is reused in local, interactive, and browser based deployments. Workarounds and pragmatic fixes were applied during development to keep the system functional across diverse environments. These corrections are written down and highlighted as spots for hardening before production deployment.

VI. MODELING

The modeling techniques to used the three different inputs: URLs, SMS, and Emails. Since these data types differ in language and structure, we created two separate systems: a Stacked Ensemble Model for URLs and SMS, and a Multi-Input Neural Network for Emails.
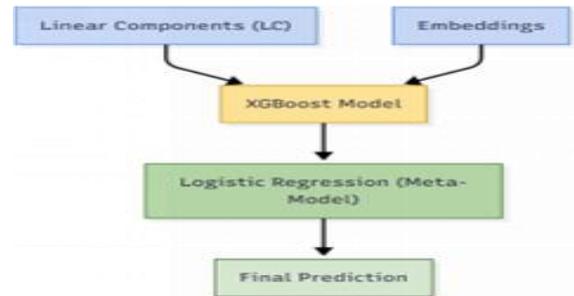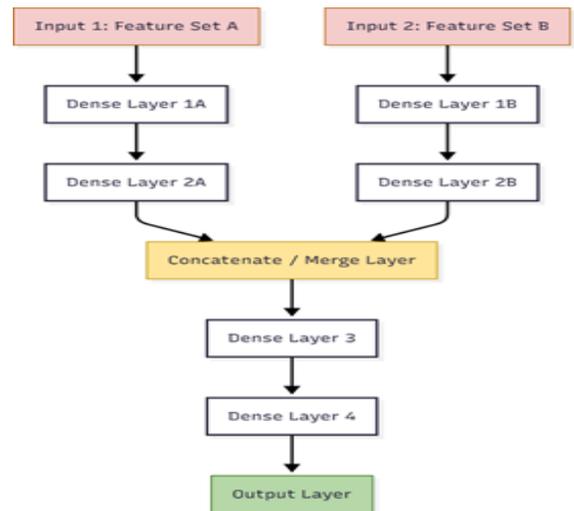


Figure 2. Stacked Ensemble Diagram



Figure 3. Multi-Input NN Architecture Diagram

URLs and SMS: Stacked Ensemble Framework: These data types have strong signals in both tabular (lexical, contextual) and semantic features. A stacked ensemble was chosen to combine these:

- Level 0 Model: For textual data such as URLs and SMS, both lexical and contextual characteristics serve as powerful discriminators. A two-tier Stacked Ensemble was developed to effectively combine handcrafted and learned features. At the base level, an XGBoost classifier was trained using lexical indicators.
- Level 0 Features: 384-dimensional transformer embeddings were used to capture semantic features, which correspond to deeper language-level patterns.
- Level 1 Model (Meta-model): These embeddings were concatenated with the probabilistic outputs of XGBoost, i.e., the embeddings were joined with the XGBoost probabilities, and a Logistic Regression model was then used to predict phishing probability. The Logistic Regression model took these concatenated inputs and gave the final phishing probability.

- Prediction: The different components of the ensemble, i.e., the symbolic and semantic information, were efficiently integrated by the ensemble through this architecture, thus it was able to achieve improved precision on heterogeneous text data.

Emails: Multi-Input Neural Network : Emails are a more complex case since they have both structured metadata and unstructured semantic text. To handle these features, a Multi-Input Neural Network (NN) was created with PyTorch. The design featured two parallel input routes:

- The outputs from the two streams were combined and fed through dense layers for joint feature learning. The last layer was a single neuron with a sigmoid activation function, optimized with LogitsLoss. A weighted loss adjustment (pos_weight) was used to handle the class imbalance, and the tabular data were normalized with a Standard Scaler. This configuration allowed the use of both language and numeric signals simultaneously, which, in turn, made the model more powerful in making subtle distinctions in the case of complex phishing emails.
- Training: Models were trained independently for each data type using the processed training and validation datasets. Trained models, including XGBoost models (.pkl), Logistic Regression meta-models (.pkl), PyTorch NN state dicts (.pt), fitted StandardScalers (.pkl), and lists of feature names (.pkl), were saved to the models/ directory.
- Challenges: One of the major aspects of the training was the handling of large datasets. The issues that arose due to class imbalance were also quite difficult to deal with, along with the selection of the appropriate hyperparameters and making sure that the data pipelines were correctly feeding features to the right models. The most difficult problem was fixing the compatibility issues when the saved PyTorch model was being loaded with the model architecture defined in code (size mismatch RuntimeError), which demanded that the model definition and saving process be in perfect agreement.

## VII. THREAT INTELLIGENCE INTEGRATION

Threat Intelligence (TI) mainly offers a list of known malicious indicators such as URLs, domains, and IPs. Using this integration provides a simple, quick, and reliable way to identify those threats that have already been detected. Thus, supporting the ML models that are more efficient in finding new or slightly changed attack cases.

Sources: We depended on publicly available data streams: OpenPhish (a list of URLs in plain text) and PhishStats (a CSV file containing URLs and scores). The URLHaus API was also considered a possible real-time lookup source.

Methodology: To regularly fetch the newest feeds from the remote sources and store them locally, we devised a separate program. The files where this program saves the data are separate from the main detection system.

To handle these local feed files that contain threat information, the shieldphish.threat_intel.ti_lookup module was created. It reads them into memory (with Python sets for quick O(1) lookup) the first time the module is imported. The perform_threat_intel_lookup function first attempts to find a match for the input (mainly URLs) locally in the feed data that is loaded into memory. This is extremely fast. If the item is not located in the feeds, a call to URLHaus API is made for a live check (with caching to allow a certain number of API calls).

In the case of emails and SMS, some dummy logic was

created just to show the way of implementing TI by getting URLs/domains from the content and then doing URL TI checks on them.

Integration into Prediction: TI lookup is the operation that takes place in the main prediction module just before model inference. If a source of phishing that is known to a phishing entity is found, the system thus goes on to return a phishing classification that is, ML processing is bypassed. Otherwise, the respective model is called for prediction. Such a dual-layer mechanism of a rule-based check followed by model-based prediction performance is optimised in that it ensures both the immediate recognition of known threats and the adaptive detection of new ones.

Challenges : The TI integration stage was very difficult and it presented the team with numerous challenges such as inconsistencies in the feed formats, API latency, and feed unavailability at certain times. The team resolved these issues in part through local caching, asynchronous update scheduling, and format unification.
The hybrid TI–ML model therefore achieved a compromise between detection speed, accuracy, and the ability to recover from phishing attacks.

## VIII. EVALUATION

The models were evaluated on unseen test sets (15% of the total processed data) after training.
- URL Model (XGBoost Ensemble): The URL model has achieved very high-performance metrics on the test set
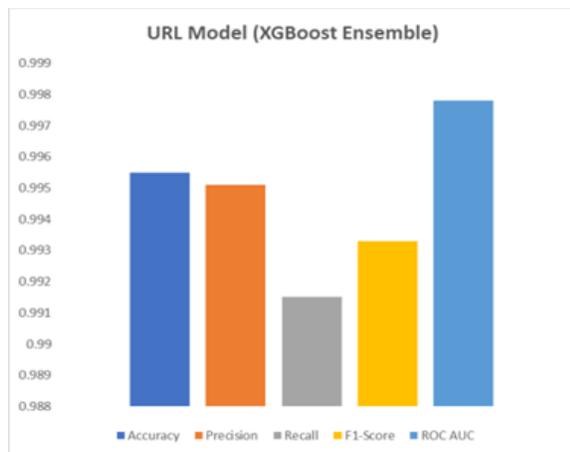


Figure 4. URl Model Performance Metrics

Discussion: Great metrics, which shows a high power of the model in differentiating the URLs. It was acknowledged that the problem of false positives for some particular real sites such as learner.vierp.in, thus, pointing to the necessity of threshold tuning and maybe more benign data variety.

- Email Model (Multi-Input NN): The model was able to deliver a good performance after refinement.



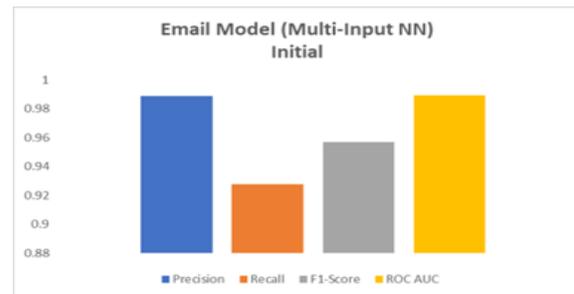Figure 5. Accuracy Improvement after refinement



Figure 6. Email Model Performance Metrics

Discussion: Initial the performance was not bad but after some changes to feature extraction (for example by using the email structure or links) and model refinement the accuracy was increased substantially. The NN's power to merge various feature types was quite helpful.
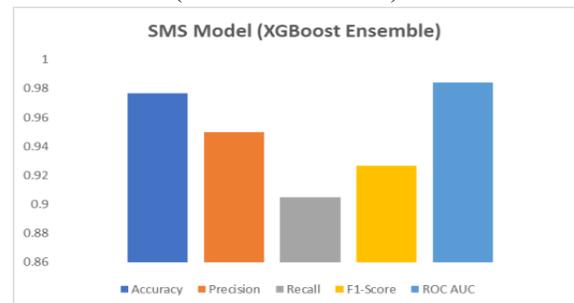
- SMS Model (XGBoost Ensemble):



Figure 7. SMS Model Performance Metrics

- XAI Evaluation: SHAP explanations for URL and SMS models give feature importance rankings.

LIME explanations for the Email model indicate the most significant words/phrases in the text. Challenges: It was observed that the behavior of the SHAP output for certain URL instances was not consistent (for example, opposing feature contributions vs actual values), thus suggesting that there might be some limitations or bugs in the way the SHAP library interacts with the model/data. The challenge of specifying the prediction function for LIME with the Multi-Input NN has been resolved.

## IX. DISCUSSION AND FUTURE WORK

ShieldPhish is an excellent instance of a working multi-modal phishing detection system that also includes integrated threat intelligence and explainability across multiple user interfaces. Despite achieving high figures on the test data, the project has led to a number of issues being opened for discussion and further research work:

- Improving Email Feature Robustness: More fine-grained parsing of the raw email MIME structure and headers can uncover more contextual features (SPF/DKIM/DMARC results, sender IP reputation, attachment analysis) that can be utilized. To fully understand the embedded components (images, complex HTML, obfuscated links) the implementation of advanced parsing methods is necessary.
- Enhanced URL Feature Engineering: Even though the current features are quite strong, the performance can still be raised by incorporating some sophisticated methods like analyzing domain generation algorithms (DGAs), verifying page content/forms (needs a crawler), or using domain/IP reputation scores from paid services to get higher accuracy and fewer false positives.
- Resolving SHAP Inconsistency: Understanding the reason of SHAP output being illogical in the case when it is used as an explainability component should be very reliable, is absolutely necessary. It might involve inspecting the code of the SHAP library, trying other SHAP explainers, or even using less complex models if they happen to be the source of the issue.
- Model Tuning and Optimization: If model hyperparameters are tuned to an even higher degree for all models and other architectures (e.g., fine-tuning transformers for all text types, more complex ensembles) are also tried, there may be very slight improvements in performance. If the Chrome Extension navigation interception is to work smoothly, making the model inference as fast as possible is very important.
- Dashboard and Logging: The Implementing a database for storing scan results together with the creation of the dashboard endpoints and Visuals would be a great medium to comprehend system usage and the changing threat landscape over time.
- API Security and Scalability: The Backend API in a real scenario will require a set of security measures like authentication, rate limiting, and it should be placed on scalable infrastructure (for example, cloud services) through production-ready WSGI servers (such as Gunicorn).
- Chrome Extension Refinement: The next steps to the work may involve improving the UI/UX, adding user configuration features (e.g., sensitivity threshold, API key), and perhaps employing more advanced browser APIs for performance or deeper page interaction. The complexity of the situation is that there are also many navigation edge cases that need to be handled.
- Continuous Learning: Concept development for the models to be retrained with newly found phishing samples (e.g., user reports) is an idea that will help the system become more secure as it will be able to adapt to the evolution of threats.

## X. CONCLUSION

ShieldPhish can be seen as an example of a possible comprehensive phishing detection system that analyses URLs, emails, and SMS messages. The project presents an effective way of protecting the digital world by developing a tool that combines feature engineering, machine-learning models tailored, threat intelligence, and explainable AI integrated into the tool, and deploying interfaces through a CLI, Streamlit app, and Chrome Extension with the ability of proactive blocking. Although substantial efforts have been made to deal with issues related to data complexity, model interpretation robustness, and environment-specific packaging, which have been addressed with pragmatic solutions, the performance metrics obtained and the functionality

across multiple platforms, thus, the system can become a powerful weapon in the Figureht against phishing.

## REFERENCES

[1] A. Sharma, S. Gupta, and R. Singh, "Detection and classification of phishing websites using machine learning: A robust framework," Journal of Cyber Security Technology, vol. 9, no. 2, pp. 87–104, May 2025.

[2] M. S. Alzboon, A. Alzain, and S. Alsmadi, "Detection of phishing websites using a machine learning algorithm," AIP Conference Proceedings, vol. 3237, p. 030059, 2024.

[3] K. Omari, "Comparative Study of Machine Learning Algorithms for Phishing Website Detection," International Journal of Advanced Computer Science and Applications, vol. 14, no. 9, pp. 356–362, 2023.

[4] T. G. Nguyen, M. A. Nguyen, and D. T. Tran, "In-Depth Analysis of Phishing Email Detection: Evaluating Machine Learning and Deep Learning Approaches," Applied Sciences, vol. 15, no. 6, p. 3396, 2025.

[5] A. Sharma, M. S. Alzboon, and S. Gupta, "Phishing Website Detection Using Machine Learning," Global Research Journal, vol. 7, no. 2, pp. 112–125, 2024.

[6] S. M. Dipto, M. T. Reza, M. N. J. Rahman, M. Z. Parvez, and P. D. Barua, "An XAI Integrated Identification System of White Blood Cell Type Using Variants of Vision Transformer," in Proceedings of the Second International Conference on Innovations in Computing Research (ICR'23), Springer, pp. 303–315, 2023.

[7] J. Zhang, Y. Wu, and L. Chen, "Transformer-based Multi-modal Phishing Detection in Emails and SMS," IEEE Access, vol. 12, pp. 45678–45689, 2024.

[8] N. Patel, R. Kumar, and A. Singh, "Explainable Transformer Models for Phishing URL Detection," IEEE Transactions on Information Forensics and Security, vol. 19, pp. 1234–1246, 2024.

[9] H. Lee, S. Kim, and J. Park, "Threat Intelligence-Driven Phishing Detection using Deep Learning," Computers & Security, vol. 135, p. 103456, 2024.

[10] M. R. Hassan, S. Rahman, and F. Islam, "Multi-modal Deep Learning for Phishing Detection across Email, SMS, and Web," Pattern Recognition Letters, vol. 177, pp. 1–8, 2024.

[11] L. Wang, Q. Li, and Y. Zhou, "BERT-Based Phishing Email Detection with Explainable AI," Expert Systems with Applications, vol. 237, p. 120456, 2024.

[12] S. Gupta, A. Sharma, and P. Kumar, "Threat Intelligence and Machine Learning for Real-Time Phishing Detection," Future Generation Computer Systems, vol. 154, pp. 98–109, 2024.

[13] R. Singh, M. Jain, and V. Sharma, "Phishing Detection in SMS using Deep Learning and NLP," IEEE Internet of Things Journal, vol. 12, no. 3, pp. 2345–2356, 2024.

[14] T. Al-Dala'in and J. H. S. Zhao, "Overview of the Benefits Deep Learning Can Provide Against Fake News, Cyberbullying and Hate Speech," in Proceedings of the Second International Conference on Innovations in Computing Research (ICR'23), Springer, pp. 13–27, 2023.

[15] P. Kumar, S. Singh, and R. Kaur, "Real-Time Phishing Detection Using Multi-Modal Transformers," IEEE Transactions on Dependable and Secure Computing, vol. 22, no. 1, pp. 45–58, 2025.

[16] Y. Liu, X. Zhang, and H. Xu, "Explainable AI for Phishing Detection: A Survey and New Directions," ACM Computing Surveys, vol. 57, no. 2, pp. 1–39, 2025.