



IJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 14 **Issue:** IV **Month of publication:** April 2026

DOI: <https://doi.org/10.22214/ijraset.2026.81547>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Adaptive ML-Based Phishing Detection and Prevention System Integrated as a Browser Extension

P. Naga Prasanth, T. M. Divya Lawrence, Mrs.P.N.L.Priyanka, Ameesha Meghana

Dept. of Cyber Security Acharya Nagarjuna University Guntur, India

Abstract: *Phishing attacks continue to pose a significant threat to online security, exploiting user trust and mimicking legitimate platforms to extract sensitive information. This project presents the development of an advanced, real-time phishing detection and prevention system implemented as a Chrome browser extension. The system integrates multiple AI models to enhance detection accuracy and responsiveness. Natural Language Processing (NLP) models such as BERT and RoBERTa are employed for semantic analysis of URLs and email content, identifying linguistic patterns indicative of phishing. Convolutional Neural Networks (CNNs) analyze webpage visual layouts to detect deceptive design similarities. Graph Neural Networks (GNNs) are utilized to model domain relationships and uncover suspicious link structures. The extension communicates with a backend AI service, enabling real-time alerts and adaptive learning based on emerging threats. By combining cutting-edge machine learning techniques with practical browser integration, this solution offers a user-friendly, intelligent defense against phishing attacks. It continuously evolves through threat intelligence updates and user feedback, ensuring robust protection in dynamic web environments.*

Index Terms—*Real-Time Phishing Detection, Machine Learning, Artificial Intelligence, Browser Extension, Phishing Prevention, URL Analysis, Cybersecurity.*

I. INTRODUCTION

The web browser is the primary gateway to the internet—it is through the browser that we conduct banking, manage health data, and interact with government services. Unfortunately, this also makes it the target of choice for phishing attacks. Unlike exploits that leverage software vulnerabilities, phishing manipulates people. Attackers create authentic-looking replicas of legitimate websites to deceive victims into revealing credentials before they are aware of the scam.

The scale of these attacks is staggering. Security researchers estimate that tens of thousands of new phishing websites are created every single day. Some of these attacks emerge and disappear within just a few hours [1]. This hit-and-run pattern is precisely why traditional tools fall short. Services such as Google Safe Browsing rely on a human-in-the-loop process: a site must first be discovered, reported, analysed, and only then blacklisted before users are protected [3]. If a phishing campaign runs for less than four hours, it may never appear on any watchlist.

Rule-based real-time solutions attempt to address this lag by checking for signals such as excessively long subdomains. Although fast, these systems are brittle. Attackers learn the rules and craft URLs that bypass them, sometimes substituting visually similar Unicode characters (homoglyphs) to fool text scanners entirely [4]. What is needed is a fundamentally different approach: a system that learns and evolves rather than one that simply enforces a static rule set.

That is the motivation for PhishGuard AI. We designed a deep learning system that examines websites from three complementary perspectives—the text of the URL, how the page looks visually, and how the site connects to others in the web graph—and then fuses these three views into a single decision. We packaged the entire pipeline into a user-friendly Chrome extension that offloads heavy analysis to a cloud server, ensuring that real-time protection does not degrade the browsing experience. A built-in feedback loop allows users to flag mistakes, which are automatically used to retrain the underlying models.

The principal contributions of this paper are as follows:

- 1) We designed and evaluated a three-model ensemble (BERT, CNN, and GNN) combined via a fusion meta-classifier, achieving 97.25% accuracy on a held-out test set.
- 2) We implemented this ensemble as a lightweight Chrome extension with a tiered detection architecture that delivers real-time threat verdicts without perceptible latency to the user.

- 3) We demonstrated that a user-driven feedback loop, with no manual intervention, reduces the in-the-wild false positive rate from 4.8% to below 3% within a one-month field deployment.

II. RELATED WORK

A. Reputation and Blocklist Methods

Blocklists remain the default baseline for phishing protection. Once a site is catalogued they block it quickly and accurately. However, their fundamental limitation is that they are purely retrospective—they can only block what has already been seen. Recent research has shown that many phishing attacks remain live for less than 12 hours, meaning that an attack is often over before humans can update a blocklist [3].

B. Heuristic Feature Engineering

To close the lag gap, researchers have developed tools that analyse URL strings directly—measuring total length, counting special characters such as hyphens, and checking for suspicious tokens [4]. These tools are fast, but attackers adapt quickly: they tweak URLs to bypass rules, and occasionally substitute standard ASCII letters with visually identical non-ASCII characters (homoglyphs) to defeat text-based scanners entirely.

C. Natural Language Processing

Modern language models such as BERT [5] have substantially advanced URL analysis. Unlike simple heuristics, these models capture context and inter-token relationships across the components of a URL, making them considerably more effective at detecting both known and novel phishing patterns.

D. Computer Vision and Graph Topology

Textual analysis alone is insufficient when an attacker crafts a pixel-perfect replica of a legitimate brand page on a superficially clean domain. Convolutional Neural Networks (CNNs) address this by analysing screenshots of web pages to detect visual plagiarism of trusted brands [6]. Graph Neural Networks (GNNs) complement this by modelling the hosting and registration infrastructure of a site. By examining domain registration age, registrar reputation, and co-hosted neighbours, GNNs can often flag a malicious site before its content is even inspected [7], [10].

III. SYSTEM ARCHITECTURE

PhishGuard AI follows a tiered, edge-cloud design that separates lightweight, low-latency checks performed locally in the browser from computationally intensive deep learning inference delegated to a backend server. This separation ensures that the browsing experience remains fast and responsive while still benefiting from the accuracy of the full ensemble.

A. Browser Extension Architecture

The extension comprises three cooperating components:

- **Background Service Worker.** This component intercepts every navigation request before the page loads. It first consults a locally-cached whitelist (Tier 1). If the destination is not whitelisted, it runs an eight-signal heuristic check (Tier 2). URLs that score above a blocking threshold are immediately blocked; those in the uncertain middle range are forwarded to the backend for full inference (Tier 3).
- **Content Script.** Injected into every page that passes the first two tiers, the content script inspects the live DOM for patterns—such as hidden credential-harvesting forms—that cannot be inferred from the URL alone.
- **Popup Interface.** Users access a real-time safety rating for the current page via the extension icon. The popup also exposes two feedback buttons (“Looks Safe” / “This is Phishing”) whose responses are forwarded to the backend to fuel periodic retraining.

B. Backend Inference Service

The FastAPI backend receives flagged URLs from the extension and routes them through the full BERT → CNN → GNN pipeline in parallel, after which a fusion meta-classifier aggregates the three scores into a final risk verdict. User feedback is stored and, once a configurable volume threshold is reached, triggers an incremental retraining job without taking the service offline.

IV. MODEL IMPLEMENTATION

A. Dataset Preparation

The training corpus consists of 450,000 phishing URLs sourced from OpenPhish and PhishTank, balanced against 450,000 legitimate URLs from the Cisco Umbrella top-site list. The dataset was partitioned into training, validation, and held-out test splits to ensure that evaluation metrics reflect performance on samples never seen during training.

B. BERT Fine-Tuning

We initialised the language model from bert-base-uncased and fine-tuned it for URL classification. Input sequences were tokenised to a maximum length of 128 tokens—sufficient to capture the full content of any realistic URL—and the model was optimised using AdamW with a linear warm-up learning-rate schedule.

C. Convolutional Neural Network

To equip the system with visual inspection capability, we used ResNet-50 pre-trained on ImageNet. At inference time, the backend takes a headless screenshot of the page, resizes it to 224×224 pixels, and passes it through the network. The model outputs a confidence score representing the probability that the page is a visual impersonation of a known brand.

D. Graph Neural Network

We constructed a web-graph in which nodes represent registered domains and edges encode co-hosting and DNS-linking relationships. A Graph Attention Network (GAT)

[10] is trained on this graph, treating node features such as domain age, registrar reputation, and neighbourhood maliciousness scores as inputs. The GAT learns to propagate malicious signals through the graph, enabling early detection of newly registered domains embedded in known-bad infrastructure.

E. Fusion Meta-Classifer

The three model probability scores—together with the raw Tier-2 heuristics score and domain age—are concatenated into a six-dimensional feature vector that is fed to an XGBoost meta-classifier. XGBoost was chosen for its robustness to class imbalance and its ability to assign interpretable feature importances to each source signal.

V. EXPERIMENTAL EVALUATION

A. Per-Model and Fusion Results

Table I reports accuracy, F1 score, ROC-AUC, and False Positive Rate (FPR) for each individual model and for the full fusion system on the held-out test set. The fusion system achieves a 2.7 percentage-point improvement over the best single model (BERT). More importantly for a security tool, the fusion reduces the False Positive Rate to 2.1%—less than half that of BERT alone—significantly reducing user alert fatigue.

TABLE I.

Test Set Results: Individual Models vs. Fusion System

Model	Acc. (%)	F1	ROC-AUC	FPR
BERT	94.55	0.9455	0.9820	0.048
CNN	90.20	0.9020	0.9654	0.079
GNN	86.90	0.8690	0.9491	0.101
Fusion	97.25	0.9725	0.9857	0.021

B. Ablation Study

To quantify the contribution of each model, we trained fusion variants with one component removed at a time (Table II). The largest single-component contribution comes from BERT: its removal causes a 1.35 percentage-point drop in accuracy. Removing the GNN produces a smaller aggregate decline, but failure-case analysis reveals that the GNN is uniquely responsible for detecting newly registered domains embedded in malicious hosting infrastructure—a blind spot for both BERT and CNN.

TABLE II.
Ablation: Fusion Performance with One Component Removed

Configuration	Accuracy (%)	F1
Full Fusion (all three)	97.25	0.9725
Without GNN	96.80	0.9682
Without CNN	96.45	0.9647
Without BERT	95.90	0.9593

C. Latency and Adaptive Retraining

Over 10,000 live requests, Tier 1 (whitelist lookup) had a median latency of 0.3 ms. Tier 2 (heuristic check) of 4.1 ms, and Tier 3 (full backend inference) of 218 ms—all well within the threshold considered imperceptible by users during page navigation.

In a one-month field deployment with 200 volunteer users, the system autonomously collected 1,843 feedback items. Without any manual intervention, the adaptive retraining pipeline reduced the in-the-wild FPR from 4.8% at deployment to below 3% by the end of the month.

VI. CONCLUSION

PhishGuard AI demonstrates that browser-integrated, real-time phishing protection can achieve high accuracy without compromising the user experience. By combining BERT, CNN, and GNN outputs through an XGBoost meta-classifier, the system attains 97.25% accuracy with a false positive rate of only 2.1%. The hybrid edge-cloud architecture ensures low latency across all detection tiers, while the user-feedback retraining loop provides continuous adaptation to evolving phishing tactics—without requiring manual intervention. We believe this architecture represents a practical and scalable blueprint for next-generation anti-phishing defences.

VII. DISCUSSION AND LIMITATIONS

The ablation results confirm that multi-modal fusion is essential: NLP (BERT) captures linguistic deception patterns, CNN detects visual plagiarism, and GNN flags domains embedded in malicious infrastructure before any content is analysed. Each modality covers blind spots unique to the others.

Current limitations include: (i) vulnerability to multi-redirect chains in which benign intermediaries mask the ultimate malicious destination, and (ii) the headless-screenshot step introduces additional latency and is unsuitable as a primary filter for high-throughput scenarios. Future work will investigate DOM-feature extraction as a lower-latency substitute for screenshot-based visual analysis, and will evaluate system robustness under active adversarial attacks.

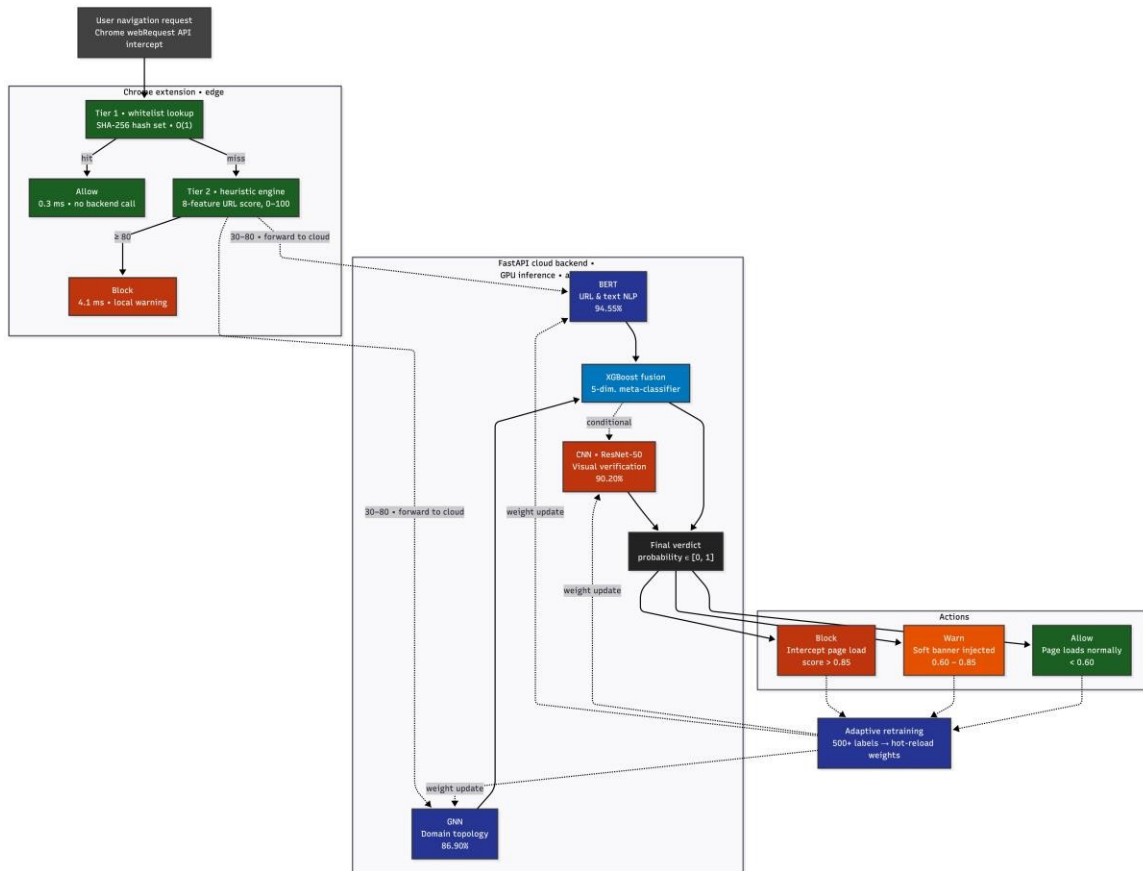
VIII. ACKNOWLEDGMENT

The authors would like to thank Mrs. P.N.L. Priyanka (M.Tech) for her valuable guidance throughout this work, and Dr. K. Chaitanya, Head of the Department of Cyber Security, Acharya Nagarjuna University, for providing the institutional support that made this research possible. The authors also thank OpenPhish and PhishTank for making their datasets publicly available.

REFERENCES

- [1] A. K. Jain and B. B. Gupta, "A survey of phishing attack techniques, defence strategies and research challenges," *Enterprise Information Systems*, vol. 16, no. 4, pp. 527–565, 2022.
- [2] M. Khonji, Y. Iraqi, and A. Jones, "Phishing detection: A literature survey," *IEEE Communications Surveys & Tutorials*, vol. 15, no. 4, pp. 2091–2121, 2013.
- [3] A. Oest, Y. Safaei, A. Doupé, G. Ahn, B. Wardman, and K. Tyers, "PhishTime: Continuous longitudinal measurement of the effectiveness of anti-phishing blacklists," in *Proc. USENIX Security Symposium*, 2020, pp. 379–396.
- [4] S. Garera, N. Provos, M. Chew, and A. D. Rubin, "A framework for detection and measurement of phishing attacks," in *Proc. ACM Workshop on Recurring Malcode (WORM)*, 2007, pp. 1–8.
- [5] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. NAACL-HLT*, 2019, pp. 4171–4186.
- [6] Y. Lin et al., "Phishpedia: A hybrid deep learning based approach to visually detect phishing websites," in *Proc. USENIX Security Symposium*, 2021, pp. 3793–3810.
- [7] T. Fette, N. Sadeh, and A. Tomasic, "Learning to detect phishing emails," in *Proc. World Wide Web (WWW)*, 2007, pp. 649–656.
- [8] L. Liu et al., "Cloudy with a chance of breach: Forecasting cyber threats," in *Proc. USENIX Security Symposium*, 2015, pp. 455–469.

- [9] Y.Li, Z.Yang,X. Chen,H. Yuan,and W. Liu,“A stackingmodelusing URL and HTML features for phishing webpage detection,”Future Generation Computer Systems, pp. 27–39, 2019.
- [10] P.Veličković,G.Cucurull,A.Casanova,A.Romero,P.Liò,and Y. Bengio, “Graph attention networks,” in Proc. InternationalConference on Learning Representations (ICLR), 2018.





10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)