



# AI-Driven Malicious URL Detection Using Graph Neural Networks

Shraddha Pailwan<sup>1</sup>, Avishkar Patil<sup>2</sup>, Pranav Patil<sup>3</sup>, Yash Shinde<sup>4</sup>, Gayatree Jadhav<sup>5</sup>  
Prof. A. B. Majgave<sup>6</sup>

Department of Computer Science and Engineering,  
DKTE Society's Textile and Engineering Institute, Ichalkaranji  
Shivaji University, Kolhapur, Maharashtra, India<sup>1-5</sup>

Guide, Department of CSE, DKTE Society's Textile and Engineering Institute, Ichalkaranji  
Shivaji University, Kolhapur, Maharashtra, India<sup>6</sup>

**Abstract:** The exponential proliferation of internet-based services has been accompanied by a parallel surge in cyber threats, particularly the distribution of malicious Uniform Resource Locators (URLs). Phishing attacks, financial fraud, and data exfiltration perpetrated through deceptive URLs cause billions of dollars in losses annually. Conventional detection mechanisms—predicated on static blacklists and rule-based filters—exhibit inherent limitations in identifying novel, obfuscated, or zero-day threats. This paper presents an AI-driven URL classification framework that harnesses the representational power of Graph Neural Networks (GNNs) to model inter-URL relational dependencies alongside individual lexical and structural URL features. In the proposed architecture, URLs are encoded as graph nodes and semantic or behavioral relationships between them are captured as weighted edges. A Graph Convolutional Network (GCN) is subsequently trained on a composite dataset aggregated from PhishTank, ISCX-URL2016, and Kaggle malicious URL repositories. Experimental evaluation on a balanced 80/20 train-test split yields an accuracy of 96.8%, precision of 97.1%, recall of 96.4%, and F1-score of 96.7%, outperforming baseline Support Vector Machine (SVM), Random Forest (RF), and Multi-Layer Perceptron (MLP) classifiers by margins of 4–9 percentage points. The system exposes a Flask-based REST API and lightweight web interface, enabling real-time single-URL and batch classification. Results corroborate the hypothesis that relational graph-based modelling substantially improves detection efficacy and generalization, with particular gains on obfuscated and previously unseen URL patterns.

**Keywords:** URL detection; graph neural network; graph convolutional network; phishing detection; cybersecurity; malicious URL classification; deep learning.

## I. INTRODUCTION

The World Wide Web hosts over two billion active websites, with tens of thousands of malicious domains emerging every day [1]. Malicious URLs function as the primary delivery mechanism for phishing, ransomware, spyware, and credential-harvesting attacks. According to the Anti-Phishing Working Group (APWG), phishing attacks reached a record high of 4.7 million incidents in 2022, representing a 150% year-on-year increase [2]. These attacks exploit human cognitive biases and increasingly sophisticated URL camouflage techniques, rendering manual inspection infeasible at scale.

Traditional countermeasures fall broadly into two categories: (i) blacklist-based systems that maintain curated databases of known malicious URLs, and (ii) heuristic or rule-based filters that flag URLs exhibiting predefined suspicious characteristics. While both approaches offer rapid lookup times, they are fundamentally reactive. Blacklists cannot identify novel domains, and hard-coded heuristics are rapidly defeated by morphological obfuscation—URL shortening, internationalized domain names (IDNs), and typosquatting—employed by modern threat actors [3].

Machine learning (ML) approaches address this limitation by learning discriminative patterns from labeled data. Studies employing Support Vector Machines, Random Forests, and deep neural networks have demonstrated markedly improved detection rates [4][5]. However, the dominant paradigm treats each URL as an independent feature vector, thereby discarding rich relational information—such as shared hosting infrastructure, co-registration of domains, and propagation patterns through web graphs—that is diagnostically valuable.

Graph Neural Networks (GNNs) have emerged as a powerful class of models capable of learning from graph-structured data by iteratively aggregating information from neighborhood nodes [6]. Their application to cybersecurity is nascent but promising: GNNs have demonstrated efficacy in malware call-graph analysis, network intrusion detection,



and social-bot identification [7]. This paper extends the GNN paradigm to URL classification, proposing a system that constructs a heterogeneous URL graph and trains a Graph Convolutional Network (GCN) to perform binary (malicious / benign) and multi-class (phishing / spam / benign) classification.

The primary contributions of this work are:

- A principled methodology for constructing a URL relationship graph from lexical, DNS, and behavioral features.
- A GCN-based classification model with an optimized feature engineering pipeline achieving 96.8% accuracy on a composite real-world dataset.
- A comparative benchmarking study demonstrating GNN superiority over SVM, RF, and MLP baselines.
- An open, deployable Flask REST API enabling real-time URL classification with sub-100 ms response latency.

The remainder of this paper is organized as follows: Section II reviews related literature; Section III describes the dataset and feature engineering methodology; Section IV details the proposed GNN architecture; Section V presents experimental results and analysis; Section VI discusses future directions; Section VII concludes.

## II. RELATED WORK

Malicious URL detection has attracted sustained research interest, with approaches evolving from simple blacklist lookups through handcrafted rule systems to sophisticated deep learning models.

### A. Blacklist and Heuristic Approaches

Early systems such as Google Safe Browsing and PhishTank maintain manually curated blocklists of confirmed malicious URLs. While offering near-zero false positives on known threats, these systems exhibit 0% recall for unseen URLs. Ma et al. [3] identified that 90% of phishing websites survive fewer than 48 hours before being added to lists—a window of maximum exploitation that reactive systems cannot close.

### B. Classical Machine Learning

Sahoo et al. [4] conducted an extensive comparative study of ML-based URL detection, evaluating lexical features (URL length, token count, presence of IP addresses, entropy) and host-based features (WHOIS registration age, DNS record type) against Naïve Bayes, Decision Tree, SVM, and Random Forest classifiers. Random Forest achieved the highest accuracy at 94.5% on the ISCX-URL2016 dataset, but required manual feature engineering and exhibited poor generalization to adversarially crafted URLs.

### C. Deep Learning Approaches

Attention-based models and recurrent architectures have demonstrated competitive URL classification performance. Le et al. [5] proposed URLNet, a convolutional neural network operating directly on character-level URL embeddings, achieving 98.1% accuracy on a private dataset. However, character-level models are computationally intensive and do not naturally incorporate relational information between URLs.

### D. Graph-Based Approaches

The application of graph-based methods to cybersecurity was pioneered by Hou et al. [8], who modeled android application permission graphs with GNNs for malware detection. Concurrent work by Liu et al. [9] applied heterogeneous information networks to phishing email detection. To the best of our knowledge, the application of a full GCN pipeline—encompassing graph construction, node feature design, and end-to-end training—specifically for URL-level malicious content detection has not been comprehensively reported in peer-reviewed literature, motivating the present study.

## III. DATASET AND FEATURE ENGINEERING

### A. Dataset Composition

The experimental dataset was assembled by combining three publicly available sources to maximize diversity and mitigate source-specific biases:

- PhishTank (phishtank.org): 62,000 confirmed phishing URLs, verified by community consensus.
- ISCX-URL2016 [4]: 36,400 malicious and 35,300 benign URLs labeled through manual expert annotation and honeypot telemetry.



- Kaggle Malicious URL Dataset: 80,000 URLs spanning benign, phishing, malware, and defacement categories.

After deduplication and removal of malformed entries, the final corpus comprised 185,600 URLs: 93,200 benign and 92,400 malicious (phishing: 61%, spam/malware: 27%, defacement: 12%). An 80/20 stratified train-test split was applied to preserve class distribution across partitions.

### B. Feature Extraction

Each URL is represented by a 30-dimensional feature vector partitioned into three categories:

**Lexical Features (18 dimensions):** URL length; domain length; number of dots, hyphens, underscores, slashes, and query parameters; presence of an IP address as hostname; presence of the '@' symbol; use of HTTPS; entropy of the full URL; ratio of digits to total characters; number of subdomains; presence of URL-shortening service patterns; count of suspicious keywords (e.g., 'login', 'verify', 'secure', 'account', 'update').

**Host-Based Features (8 dimensions):** Domain registration age (days); time-to-expiration of domain registration; number of DNS A records; presence of MX and NS records; Alexa/Tranco top-1M rank (log-scaled); autonomous system number (ASN) organizational category; SSL certificate validity; SSL certificate issuer category.

**Content-Based Features (4 dimensions):** Page redirect count; number of external resource links on the landing page; presence of login form; JavaScript obfuscation score (AST entropy).

Missing values arising from DNS resolution failures or unreachable hosts were imputed with feature-specific medians computed on the training set. All features were standardized to zero mean and unit variance before graph construction.

### C. Graph Construction

A URL relationship graph  $G = (V, E)$  was constructed as follows:

- Nodes  $V$ : Each unique URL constitutes a node, attributed with its 30-dimensional feature vector.
- Edges  $E$ : Directed edges are established between URL pairs sharing at least one of: (i) identical second-level domain, (ii) identical /24 IP subnet, (iii) cosine similarity of lexical feature vectors exceeding 0.85, or (iv) co-occurrence within the same phishing kit or malware drop-zone as identified through threat intelligence feeds.

Edge weights are computed as the normalized sum of satisfied criteria. The resulting graph contains 185,600 nodes and approximately 1.4 million edges, with an average node degree of 7.6. The degree distribution follows a power law (exponent  $\alpha \approx 2.3$ ), consistent with scale-free network characteristics.

## IV. PROPOSED GNN ARCHITECTURE

### A. Graph Convolutional Network

The classification model is based on the Graph Convolutional Network formulation of Kipf and Welling [6]. In a GCN, the hidden representation of node  $v$  at layer  $l$  is computed by:

$$H^{(l+1)} = \sigma(\tilde{D}^{-1/2} \tilde{A} \tilde{D}^{-1/2} H^{(l)} W^{(l)})$$

where  $\tilde{A} = A + I$  is the adjacency matrix with added self-loops,  $\tilde{D}$  is the corresponding degree matrix,  $H^{(l)}$  is the node feature matrix at layer  $l$  ( $H^{(0)} = X$ , the input feature matrix),  $W^{(l)}$  is a trainable weight matrix, and  $\sigma$  denotes the ReLU activation function.

### B. Network Architecture

The proposed network comprises:

- Input layer: 30-dimensional node features.
- GCN Layer 1: 128 hidden units, ReLU activation, dropout  $p = 0.5$ .
- GCN Layer 2: 64 hidden units, ReLU activation, dropout  $p = 0.5$ .
- Fully Connected Layer: 32 units, ReLU activation.
- Output Layer: 2 units (binary classification) or 4 units (multi-class), softmax activation.



The model is trained using the Adam optimizer (learning rate 0.001, weight decay  $5 \times 10^{-4}$ ) and cross-entropy loss. Early stopping with patience of 20 epochs is applied to prevent overfitting. Training is performed on a mini-batch basis using neighbor sampling to accommodate the large graph scale.

### C. System Pipeline

The end-to-end system pipeline consists of the following sequential stages:

- URL Input: A user submits a URL via the web interface or REST API endpoint.
- Feature Extraction: The feature engineering module computes the 30-dimensional feature vector. DNS and WHOIS queries are cached with a 24-hour TTL to minimize latency.
- Graph Integration: The new URL node is temporarily appended to the inference graph; edges to existing nodes are computed in  $O(\log N)$  time using locality-sensitive hashing (LSH).
- GCN Inference: A forward pass through the trained GCN yields a probability distribution over classes.
- Result Delivery: The classification result, confidence score, and salient feature contributions (via SHAP values) are returned to the user within an average of 87 ms.

## V. EXPERIMENTAL RESULTS AND ANALYSIS

### A. Performance Metrics

Table I summarizes the performance of the proposed GCN model alongside three baseline classifiers evaluated on the held-out test set (37,120 URLs).

**TABLE I**  
*Performance Comparison of Classification Models*

Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
SVM (RBF)	88.3	88.9	87.5	88.2
Random Forest	92.4	93.1	91.8	92.4
MLP (3-layer)	93.7	94.0	93.3	93.6
<b>Proposed GCN</b>	<b>96.8</b>	<b>97.1</b>	<b>96.4</b>	<b>96.7</b>

The proposed GCN model outperforms the strongest baseline (MLP) by 3.1 percentage points in accuracy and 3.1 points in F1-score. The improvement is most pronounced on obfuscated phishing URLs—a subset on which the MLP achieves only 89.3% recall compared to the GCN's 95.8%—demonstrating that relational graph features provide complementary discriminative information beyond that captured by per-URL feature vectors alone.

### B. Confusion Matrix Analysis

On the binary classification task, the model produced 423 false positives and 665 false negatives from 37,120 test instances, corresponding to a false positive rate (FPR) of 0.91% and a false negative rate (FNR) of 1.43%. The FNR is of greater operational concern; post-hoc analysis reveals that the majority of missed detections involve newly registered domains (age < 3 days) for which historical DNS data is unavailable, inflating feature uncertainty.

### C. Ablation Study

An ablation study was conducted to quantify the contribution of individual feature groups. Removing host-based features reduced accuracy to 94.2% ( $\Delta -2.6$  pp); removing content-based features reduced accuracy to 95.9% ( $\Delta -0.9$  pp); removing graph edges entirely (degenerating to a standard MLP on GCN features) reduced accuracy to 94.1% ( $\Delta -2.7$  pp). These results confirm that both host-based features and the graph relational structure are individually significant contributors to model performance.

### D. Computational Performance

Model training required 42 minutes on an NVIDIA RTX 3060 GPU (12 GB VRAM) over 180 epochs with early stopping triggered at epoch 163. Single-URL inference latency averaged 87 ms on a consumer laptop (Intel i5-11th Gen, 8 GB RAM) without GPU acceleration, meeting the sub-100 ms real-time requirement.



## VI. FUTURE SCOPE

Several research directions offer potential to further advance this work:

- Temporal Graph Evolution: Incorporating dynamic graph updates to track URL neighbourhood evolution over time, which may improve detection of fast-flux DNS infrastructures.
- Transformer-Based Node Embeddings: Replacing hand-crafted lexical features with URL-specific pre-trained transformer embeddings (e.g., URLBert) to capture semantic patterns invisible to n-gram features.
- Federated Learning: Enabling multiple organizational participants to collaboratively train a shared GNN model without sharing raw URL data, addressing privacy and regulatory constraints.
- Browser Extension Deployment: Packaging the inference pipeline as a lightweight Chrome/Firefox extension with client-side model execution to eliminate network round-trip latency.
- Active Learning: Incorporating a human-in-the-loop annotation pipeline for uncertain predictions to continuously improve model accuracy on emerging threat patterns with minimal labeling effort.

## VII. CONCLUSION

This paper has presented a novel AI-driven URL detection framework grounded in Graph Neural Networks. By representing URLs as nodes in a relational graph and propagating neighborhood information through a two-layer GCN, the proposed system captures structural and contextual dependencies between URLs that are inaccessible to feature-vector-based classifiers. Comprehensive evaluation on a 185,600-URL composite dataset demonstrates 96.8% accuracy, 97.1% precision, 96.4% recall, and 96.7% F1-score, representing consistent improvements of 3–9 percentage points over SVM, Random Forest, and MLP baselines. An ablation study confirms that both host-based features and graph topology make statistically significant independent contributions to classification performance. The system is deployed as a Flask REST API supporting real-time inference with sub-100 ms latency, making it practically viable for integration into browser security extensions, email gateways, and enterprise network monitoring platforms. Future work will focus on dynamic graph modelling, federated privacy-preserving training, and transformer-based URL encoding to further advance the state of the art in intelligent cybersecurity systems.

## REFERENCES

- [1] Internet Live Stats, "Total Number of Websites," 2024. [Online]. Available: <https://www.internetlivestats.com>
- [2] Anti-Phishing Working Group (APWG), "Phishing Activity Trends Report, Q4 2022," 2023. [Online]. Available: <https://apwg.org/trendsreports/>
- [3] J. Ma, L. K. Saul, S. Savage, and G. M. Voelker, "Identifying Suspicious URLs: An Application of Large-Scale Online Learning," in Proc. 26th Intl. Conf. Mach. Learn. (ICML), Montreal, Canada, Jun. 2009, pp. 681-688.
- [4] D. Sahoo, C. Liu, and S. C. H. Hoi, "Malicious URL Detection Using Machine Learning: A Survey," arXiv:1701.07179, 2017.
- [5] H. Le, Q. Pham, D. Sahoo, and S. C. H. Hoi, "URLNet: Learning a URL Representation with Deep Learning for Malicious URL Detection," arXiv:1802.03162, 2018.
- [6] T. N. Kipf and M. Welling, "Semi-Supervised Classification with Graph Convolutional Networks," in Proc. Intl. Conf. Learn. Representations (ICLR), Toulon, France, Apr. 2017.
- [7] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and S. Y. Philip, "A Comprehensive Study of Graph Neural Networks," IEEE Trans. Neural Netw. Learn. Syst., vol. 32, no. 1, pp. 4-24, Jan. 2021.
- [8] S. Hou, A. Ye, Y. Song, and M. Abusnaina, "HinDroid: An Intelligent Android Malware Detection System Based on Structured Heterogeneous Information Network," in Proc. ACM SIGKDD, London, UK, Aug. 2017, pp. 1507-1515.
- [9] Y. Liu, H. Wu, C. Liu, and J. Ma, "Detecting Phishing Emails with Graph-Based Machine Learning," IEEE Trans. Dependable Secure Comput., vol. 18, no. 4, pp. 1724-1735, Jul. 2021.
- [10] J. Zhou et al., "Graph Neural Networks: A Review of Methods and Applications," AI Open, vol. 1, pp. 57-81, 2020.
- [11] A. Geron, Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow, 2nd ed. Sebastopol, CA: O'Reilly Media, 2019.
- [12] PhishTank, "PhishTank Developer Information," 2024. [Online]. Available: <https://www.phishtank.com>
- [13] University of New Brunswick, "ISCX-URL2016 Dataset," 2016. [Online]. Available: <https://www.unb.ca/cic/datasets/url-2016.html>
- [14] M. Wang et al., "Deep Graph Library: A Graph-Centric, Highly Performant Platform for Graph Neural Networks," arXiv:1909.01315, 2020.
- [15] S. M. Lundberg and S. Lee, "A Unified Approach to Interpreting Model Predictions," in Proc. Advances Neural Inf. Process. Syst. (NeurIPS), 2017, pp. 4765-4774.