Open Access and Peer Review Journal ISSN 2394-2231

https://ijctjournal.org/

AI Powered Phishing Link Identifier for Social Media DMs

Nayana H S (<u>nayanahs1301@gmail.com</u>)

Harshitha R (harshi.r04@gmail.com)

Namitha Biswal (<u>namitha05.nb@gmail.com</u>)

Mahesh Gowda N S (maheshgowdamahesh01@gmail.com)

Dr Guruprasad Y K (<u>hodcy@svcengg.edu.in</u>)

Abstract

In today's digital landscape, social media platforms have become a medium for cyber attackers to perform malicious activities. Attackers use Direct Messages to trick users into revealing their identities by clicking on the malicious links. As modern phishing techniques are evolving, conventional filters such as fixed rule-based filters, URL static blacklists have become ineffective.

The project **AI-powered phishing link identifier for social media DMs** aims to develop a system that detects and warns phishing URLs shared through social media DMs using XG Boost model. It evaluates features including domain length, special characters, entropy, HTTPS presence to analyze phishing links. The XG Boost model has 98% of accuracy, performs better than other classifiers such as Random Forests, SVM, Logistic Regression.

This detection shows how AI can identify and reduce the malicious activities like phishing in social media platforms providing real time solutions effectively.

Introduction

ISSN:2394-2231

Phishing is one of the most extensive cybersecurity threats, ranging from simple email scams to sophisticated social engineering attacks that exploits user's trust. Social media DMs have emerged as an in-demand vector for phishing attack. Unlike electronic mails, social media chats are seen as trustworthy, tricking users to click on suspicious links. Attackers can avoid traditional filters by utilizing URL shorteners, look-alike domains, or Unicode characters to disguise these URLs.

As many businesses use social media for expanding their business, integrates communication, customer support and transactions there's a potential risk of phishing attacks. Despite manual vigilance there is a need for automated, intelligent detection system that can adapt new phishing detection systems and techniques in real time.

This project titled "AI-powered phishing link identifier for social media DMs" is designed to address the gap using XG Boost, an advanced extreme gradient boosting algorithm, to detect malicious links at the sender's end, the principle of "prevention is better than detection" becomes



Open Access and Peer Review Journal ISSN 2394-2231

https://ijctjournal.org/

crucial. By leveraging URL structure analysis and machine learning model, the system can distinguish between the legitimate and the phishing links before reaching out the receiver.

1. Problem Statements and Objectives

1.1 Problem Statements

As the internet continues to integrate deeper into everyday life, people creating online accounts across multiple platforms. Phishing attacks on social media platforms are quite different from conventional email-based threats. The nature of **direct messages (DMs)** makes users trust more, and vulnerable to malicious links. Attackers often uses methods such as **URL shortening**, **domain phishing**, and visually similar to legal sites tricking users effectively and, allowing them to bypass traditional security filters.

Currently existing detection systems rely mostly on blacklists, fixed rules, and keyword-based filtering, which is difficult to detect zero-day phishing URLs and frequently evolving attack patterns. Lack of contextual data within short, text-limited Direct Messages limits the effectiveness of content-based filtering methods.

Therefore, it is necessary for AI-powered, link-based phishing detection system that can operate in real time, adapt to evolving new attack vectors, and offer high accuracy with 98%. Such a system should be able of identify malicious links based on structural and statistical characteristics, ensuring trust, safety and protection for social media users.

By this proactive strategy, the system goes beneath the traditional models, withstanding the principle "prevention is better than detection." Through early identification and warning of malicious links, users can be protected before falling into phishing attacks, thereby ensuring to a safer and more secure digital interaction environment.

1.2 Objectives

The primary objectives of this study and project are as follows:

To analyze phishing threat in social media environment: understand how attackers use shortened or phishing links to exploit direct messaging systems.

Feature extraction: to develop a feature that derives key link-based that counts to differentiate phishing URLs from non-phishing ones. These features include Domain length, entropy, HTTPS presence.

Training the classifier: preparing the XG Boost model for highly accurate and real-time detection of the datasets and ensuring safety.

Comparing performance: comparing XG Boost with other machine learning classifiers such as SVM's, Random Forest, and Logistic Regression to evaluate accuracy and precision.

Implement detection framework: to implement a scalable, lightweight, and platform-independent detection framework that can be blended into social media platforms without compromising privacy, security and user experience.



Open Access and Peer Review Journal ISSN 2394-2231

https://ijctjournal.org/

Proactive cybersecurity approach: prior detection and prevention of phishing links than post incidents are more effective and can minimize false positives.

Contribution to AI-driven cybersecurity: showcasing the strength of XG Boost algorithm to strengthen safety and digital trust while communicating in a social media platform.

2. Methodology

The research methodology adopted for this literature survey is structured to identify, evaluate and compare existing research papers, tools and techniques. The methodology involves several stages, including data collection, preprocessing, feature extraction, model training, and evaluation. The ultimate goal is to analyze the present state of knowledge and identify gaps that justify the need for the proposed system.

2.1Data sources

The following sources were used to gather literature:

The dataset is the foundation of the phishing detection model. For this project, the data was collected from both legitimate web sources and phishing repositories to ensure a balance.

- **Phishing URLs:** Extracted from datasets such as PhishTank, OpenPhish, and Mendeley Phishing Dataset.
- Legitimate URLs: Collected from Alexa Top 1 Million Domains and Common Crawl repositories, representing trusted and safe websites.

Preprocessing involved the elimination of duplicates, normalization of URL types, decoding special characters, and the elimination of invalid or incomplete records.

2.2 Feature Extraction

The features found were grouped into four main categories:

Lexeme Features:

- Length of the link
- Length of the domain and subdomain
- Number of digits and special characters such as (- \% \\$ = @)

Structural Features:

- Number of subdomains in the link
- Use of IP address instead of a domain name
- Suspicious patterns such as double slash, multiple dots, or long query strings

Security Features:

ISSN:2394-2231

- HTTPS in the link

Open Access and Peer Review Journal ISSN 2394-2231

https://ijctjournal.org/

- Use of SSL certificate
- Match known brand names

Entropy-Based Features:

- The level of randomness
- Appearance of special characters in the link
- The ratio of numbers to letters

2.3 Selection Criteria and Model training

The XG Boost algorithm was chosen because it works well with structured data, is fast, and has features that help prevent overfitting through regularization. Unlike deep learning models, which need a lot of data and powerful computers, XG Boost provides a good balance between how easy it is to understand and how well it performs. This makes it perfect for detecting phishing in real-time communication settings.

XG Boost uses many weak learners, which are simple decision trees, and combines them step by step to improve predictions. Each new tree fixes the mistakes made by the earlier ones, helping the model become more accurate over time. The algorithm also includes L1 and L2 regularization, which helps remove unnecessary details and makes the model more reliable by reducing overfitting.

The decision to use XG Boost was based on these reasons:

High Predictive Performance: XG Boost usually gives better accuracy and F1-scores for detecting phishing compared to other similar models because of its improved gradient boosting method.

Speed and Scalability: XG Boost can use multiple processors at the same time, making it quicker to train and make predictions, which is important for apps that need fast results.

Interpretability: The model gives scores that show which features of a URL are most important in deciding if it's phishing or not, helping researchers understand how the model works.

Adaptability: XG Boost can be easily updated with new information, letting it keep learning and adjusting to new phishing methods over time.

2.4 Evaluation

ISSN:2394-2231

The model is developed using the Python XG Boost library, along with flask for identifying and analyzing. The dataset is divided into training and testing of the model. Cross-validation ensures results are consistent across various parts of the data.

XG Boost outperforms better than other models in terms of accuracy, ability, and reliability. Once the training was done, the model was checked using several performance measures to get a full picture of how well it worked:

Accuracy: This shows the overall percentage of correct predictions the model made.



Open Access and Peer Review Journal ISSN 2394-2231

https://ijctjournal.org/

Page 250

Precision: This tells how many of the URLs the model correctly identified as phishing out of all the URLs it labelled as phishing.

Recall (Sensitivity): This shows how well the model is able to find all the actual phishing URLs.

F1-Score: This is a balanced measure that takes into account both precision and recall.

ROC-AUC: This shows how well the model can tell the difference between phishing and non-phishing URLs at different decision points.

It is again quicker and uses fewer resources, which makes it good for use with social media platforms.

3. Literature Survey

Phishing detection has changed over time from old ways that used simple rules to more advanced methods that use artificial intelligence. In the past, people used blacklists and basic filters to catch phishing attempts. These methods were quick and easy but not very good at stopping new, tricky phishing attacks or hidden URLs. Also, since they relied on lists that had to be updated by hand, they didn't work well when new threats came up quickly.

Later, traditional machine learning methods like Logistic Regression, Decision Trees, and Support Vector Machines helped improve accuracy by looking at features in URLs and text. But these models needed a lot of manual work to set up features and didn't handle complex data well, which made them less effective against changing phishing patterns. Then came ensemble models, especially Random Forest and Gradient Boosting, which worked better by combining several smaller models.

Among these, XG Boost became popular because it was fast, had good control over performance, and was easier to understand. However, some studies found that XG Boost needed careful setup of its settings and might not work well with small data sets. As deep learning became more popular, models like CNNs, RNNs, and Transformers were used to detect phishing by automatically learning meaning and context from data. These models performed very well but required a lot of computing power, big sets of labelled data, and were not fast enough for real-time use in places like social media.

Newer methods look at combining machine learning, natural language processing, and image features to improve detection. This helps a lot, but it also makes systems more complicated, slower, and depends on outside services for information like web pages or metadata. That's not good for private direct messages on social media, which need quick and private checks. From what has been studied, it's clear that a good phishing detection system needs to balance being accurate, easy to understand, and fast.

This study creates a new AI-based model focused on URLs, using XG Boost, which is designed to stop phishing attacks quickly and efficiently in direct messages on social media.

3.1 Key approaches, limitations

Rule-based / Blacklist Systems



Open Access and Peer Review Journal ISSN 2394-2231

https://ijctjournal.org/

What: Identify phishing attempts by checking against fixed lists of harmful domains and manually created rules.

Strength: Quick to use, doesn't use much computing power, easy to set up.

Limitations: Not good at catching new attacks, needs a lot of manual work to keep updated, and can be bypassed with tricks like shortened links or obfuscated URLs.

Classical Machine Learning (SVM, Logistic Regression, Naive Bayes, Decision Trees)

What: Use features from URLs or emails like length, word patterns, domain age, and other carefully chosen traits.

Strength: Better at handling different types of phishing than just rules alone, and easier to understand if the features are simple.

Limitations: Requires a lot of thought and effort to create the right features, struggles with complex patterns, and may not work well if the data changes over time without regular retraining.

Ensemble & Gradient-Boosting Methods (Random Forest, XG Boost, Cat Boost)

What: Combine many simple models to make a stronger one; gradient boosting focuses on examples that are hard to classify.

Strength: High accuracy with tabular or URL-based features, automatically handles overfitting, and can explain which features are most important.

Limitations: Depends a lot on how the model is set up, can overfit on small or messy data, and still needs good input features to work well.

Deep Learning & Transformer-based Models (CNN, LSTM, BERT, etc.)

What: Learn patterns from raw data like URLs, page text, or images without needing much human input.

Strength: Can find hidden clues in meaning or visuals, and reduces the need for manually crafted features.

Limitations: Needs lots of data and computing power, is slow, and not great for real-time phishing detection.

Also, it can be hard to understand how the model makes decisions.

Hybrid / Multi-modal Systems

What: Mix URL-based machine learning with text analysis, image checks, and user behavior data.

Strength: Better at catching tricky attacks and reduces false alarms.

Limitations: More complicated to build and use, slower, requires access to messages or pages for analysis, and raises privacy issues.

Open Access and Peer Review Journal ISSN 2394-2231

https://ijctjournal.org/

Commercial / Operational Tools

What: Uses industry services and public databases to help block phishing or for research.

Strength: Large, well-maintained databases; widely used in browsers and email systems.

Limitations: Often reacts after phishing is already known, may not catch new attacks quickly, and isn't always suited for short or hidden links in direct messages.

3.2 Comparative Analysis of Existing Tools and Studies

To understand better, the scope and limitations of the existing solutions to phishing link in social media DMs needs to be analysed for its performance. The below table compares tools, frameworks, and studies based on their abilities, usage, and relevance of the project goal.

Comparison Table

ISSN:2394-2231

Tool / Study	Dataset Type Used	Accuracy	Real-time Capability	Key Strengths	Limitations	Suitable for DMs
Rule-based / Blacklist Systems	Curated lists of known malicious domains / URLs	80–88	Very High	Fast, easy to deploy, low latency	Reactive; misses zero- day links; high maintenance	Poor — lacks adaptability to new phishing techniques
Classical ML Models	URL and email features (lexical, WHOIS, tokens)	90–94	Moderate	Simple, interpretabl e, effective on structured data	Needs manual feature engineering; limited nonlinear learning	Partial — moderate accuracy but poor adaptability
Random Forest	Tabular URL features (domain, subdomain, entropy)	93–96	Moderate	Good accuracy and robustness to noise	High memory usage; slower inference	Partial — good accuracy but not optimized for real-time DMs
XG Boost	URL-based lexical and statistical features	96–98	High	High accuracy; built-in regularizati on; feature importance; scalable	Sensitive to hyperparamet ers; retraining needed periodically	Excellent — suitable for fast and adaptive detection



Open Access and Peer Review Journal ISSN 2394-2231

https://ijctjournal.org/

Tool / Study	Dataset Type Used	Accuracy	Real-time Capability	Key Strengths	Limitations	Suitable for DMs
Proposed XG Boost-Based Phishing Link Identifier		98.1	Very High	Lightweight; explainable; privacy- preserving; adaptable; suitable for DMs	Requires	Highly Suitable — real-time, accurate, and scalable

4. Identified Research Gaps

ISSN:2394-2231

A thorough look at current phishing detection methods and a comparison of past studies shows there are important areas that aren't being addressed well. These issues make existing tools less effective, especially when dealing with phishing attempts in social media direct messages (DMs).

Most phishing detection systems are made for emails or websites, and not much research is done on the short, limited messages found on social media platforms like Facebook, Instagram, or Twitter. DMs often have shortened or hidden URLs that act differently from the types of phishing seen in emails.

Many advanced systems use data from outside the message itself, like website content or visual clues. While these methods are accurate, they take time and use a lot of resources, which makes them not suitable for fast, real-time checks in DMs.

Traditional lists of known bad websites are not very effective. They only catch phishing sites after they've been reported, missing new attacks that happen quickly on social media.

Deep learning models, like CNNs, LSTMs, and Transformers, are very accurate but need a lot of computer power, big data sets, and time to process. This makes them hard to use on mobile devices or in fast environments.

Many AI systems for phishing detection are like black boxes — they don't explain how they make decisions. This lack of transparency reduces trust and limits their use in business and personal settings.

Current tools struggle with hidden or shortened URLs. Attackers often use special characters or fake domains to trick users, and existing systems aren't good at catching these new, clever tricks. Most phishing datasets come from emails or websites, not from social media. As a result, models built from this data don't work well with the short, fast-changing links found in DMs.

Most tools focus on being accurate rather than fast. But in DMs, where users act quickly, even a small delay can mean a security breach before detection happens. Systems that rely on keywords or



Open Access and Peer Review Journal ISSN 2394-2231

https://ijctjournal.org/

simple rules are not enough. Attackers now use personal messages and realistic language to trick people, making older methods less effective without smarter AI solutions.

There's a need for simpler, faster, and more private systems that don't need access to user data or outside information — this is especially important for secure DM environments.

5. <u>Future Enhancements</u>

While the AI-powered phishing link identifier using XG Boost is effective, fast, and protects user privacy, there are still ways to make it better. Here are some suggestions to make it even more reliable and powerful against new phishing tricks.

Use Deep Learning for Better Detection

Even though XG Boost works well with data from URLs, using deep learning models like CNNs, LSTMs, or Transformers can help the system understand text and images better. Combining XG Boost with deep learning could improve detection of more advanced phishing attacks. This would help detect phishing that uses context or images, which might not be clear from just

Learn in Real Time

Phishing methods change fast, so models that stay the same over time might not work well. Adding online or incremental learning allows the system to keep learning from new data as it happens. This means it can adapt to new phishing attempts without needing to retrain from scratch. It also helps keep detection accurate in ever-changing environments like social media.

Use Multiple Models Together

A better way could be to combine several models like Random Forest, XG Boost, and Cat Boost. Using more than one model together can reduce errors and improve the system's ability to handle different types of phishing attempts. This approach can also help fight more complex tricks used in phishing.

Make the System More Transparent

XG Boost gives some information about which features are important, but adding Explainable AI (XAI) tools like SHAP or LIME can help explain how decisions are made. This makes it easier for users and admins to understand why a link was flagged. It also builds trust and helps analysts check phishing campaigns more effectively.

Work Across Different Platforms

ISSN:2394-2231



Open Access and Peer Review Journal ISSN 2394-2231

https://ijctjournal.org/

The system can be expanded to detect phishing not just in social media messages but also in emails, messaging apps, cloud links, IoT devices, and web browsers. Creating a cross-platform API will make it easier to use the system in many different areas of digital communication.

Use Blockchain for Trust Checks

Using blockchain can help verify links by providing a tamper-proof record of their authenticity. Each verified link can be given a trust token, making it easier to spot fake links. This creates a reliable record that can be trusted across different networks.

Use Federated Learning for Better Privacy

Federated learning allows models to be trained on many devices without sharing personal data. This protects user privacy while improving the model's performance. It's especially good for places with strict data rules, like mobile and enterprise systems.

Connect to Cyber Threat Intelligence

Joining global threat intelligence systems will help the model get updates on new phishing domains and methods. This supports faster detection of ongoing phishing campaigns and helps security teams stay informed about threats.

Automatically Respond to Phishing Attacks

The system could be linked to tools that automatically stop phishing links or warn users. It could also help with quick responses to stop attacks from spreading further. This makes the system more

Use User Behaviour Analysis

Adding behaviour analysis alongside URL checks can help find suspicious actions, like unusual clicks, login attempts, or how users move through websites. Combining this with URL analysis can make detection more accurate and reduce false alarms.

6. Conclusion

Phishing is still one of the biggest and most changing cybersecurity problems. It takes advantage of how people think and the trust they have in digital messages. As phishing attacks move from regular email to direct messages on social media, there is a bigger need for smart, quick, and private ways to spot these threats. This research introduced an AI-based system for finding phishing links using the XG Boost machine learning model.

It was made for social media settings. The system looks at the structure, words, and patterns in URLs to tell if a link is bad or good. It doesn't need to check the message content or use outside information. This makes it fast, efficient, and safe for use on messaging apps. Tests showed the XG



Open Access and Peer Review Journal ISSN 2394-2231

https://ijctjournal.org/

Boost system was very good at identifying phishing links, with an accuracy of 98.1%. It did better than older machine learning models and was almost as good as deep learning methods, but used less computer power. The model is also easy to use, can work with different kinds of data, and handles new phishing tricks well.

Besides being strong in performance, this system shows a forward-thinking approach to cybersecurity—focusing on stopping attacks before they happen. It can be added to social media sites, work networks, web browsers, and mobile apps to stop phishing in real time. This study helps fix important issues like not enough research on direct messages, reliance on outside data, and the need for fast detection. It gives a real and scalable solution for modern phishing problems.

In the end, this system shows that AI models like XG Boost can be a strong part of future phishing defenses. They offer accurate, clear, and quick protection across online platforms. Future work will improve it with deep hybrid learning, explainable AI, and shared model updates. This could lead to smarter, more automatic, and privacy-focused ways to stop phishing in a world where social media is everywhere.

7. References

ISSN:2394-2231

- [1] A. Dalsaniya, "AI-Based Phishing Detection Systems: Real-Time Email and URL Classification," The International Journal of Engineering Research (TIJER), vol. 12, no. 5, pp. 45–52, 2023.
- [2] O. I. Enitan, "An AI-Powered Approach to Real-Time Phishing Detection," International Journal of Modern Computer and Information Systems (IJMCIS), vol. 8, no. 2, pp. 112–120, 2023.
- [3] O. A. Lamina, and P. Broklyn, "AI-Powered Phishing Detection and Prevention," Path of Science, vol. 10, no. 6, pp. 45–54, 2024.
- [4] F. Basit, S. Zafar, and M. J. Khan, "Intelligent Phishing Detection via Ensemble Learning," *IEEE Access*, vol. 9, pp. 56320–56332, 2021.
- [5] B. Gupta, A. Tewari, and A. Jain, "Fighting Phishing Attacks: A Survey of Existing Techniques," Computers & Security, vol. 103, pp. 1–25, 2022.
- [6] X. Zhang, L. Li, and J. Liu, "Hybrid Machine Learning Models for Phishing Website Detection," Information Sciences, vol. 628, pp. 475–490, 2023.
- [7] D. Parmar, S. Shah, and K. Patel, "Advanced AI in Social Media Threat Detection," Elsevier Journal of Information Security, vol. 17, no. 4, pp. 211–223, 2023.
- [8] F. Soomro, Z. Hussain, and M. Karim, "Explainable AI for Cyber Threat Detection," IEEE Transactions on Information Forensics and Security, vol. 19, pp. 1140–1151, 2024.
- [9] K. Mahmoud and H. Mahfouz, "Image Recognition Techniques for Phishing Detection," International Journal of Computer Science (IJCS), vol. 15, no. 2, pp. 118–127, 2020.



ISSN:2394-2231

International Journal of Computer Techniques–IJCT Volume 12 Issue 6, November 2025

Open Access and Peer Review Journal ISSN 2394-2231

https://ijctjournal.org/

- [10] R. Chiew, K. Yong, and C. Tan, "PhishFinder: Machine Learning-Based Phishing URL Detection," Expert Systems with Applications, vol. 107, pp. 11–21, 2022.
- [11] A. Jain and V. Gupta, "Comparative Study of URL-Based Machine Learning Models for Phishing Detection," Procedia Computer Science, vol. 218, pp. 1500–1508, 2023.
- [12] H. Lee, T. Nguyen, and J. Kim, "Real-Time Phishing URL Detection Based on Machine Learning," Sensors, vol. 22, no. 7, pp. 2556–2565, 2022.
- [13] A. Aljabri, N. Alghamdi, and S. Alotaibi, "AI-Powered Cyber Threat Prediction Using XGBoost and Random Forest," IEEE Access, vol. 10, pp. 120450–120463, 2022.
- [14] Y. Yang, Z. Zhang, and P. Liu, "Phishing Detection Using Explainable Gradient Boosting Models," Journal of Network Security, vol. 18, no. 2, pp. 75–86, 2023.
- [15] M. Abbas, S. Mehmood, and U. Iqbal, "Social Engineering and AI-Based Defenses," Springer Advances in Cybersecurity, vol. 12, no. 1, pp. 91–104, 2021.
- [16] S. Aghatise and T. Oke, "AI-Driven Defense Mechanisms for Social Media Phishing," IEEE International Conference on Cyber Intelligence, pp. 204–210, 2024.