

# Detecting Phishing Websites using recent Techniques: A Systematic Literature Review

K Subashini<sup>1\*</sup>, V Narmatha<sup>2</sup>

<sup>1</sup>Department of Computer and Information Science, Faculty of Science, Annamalai University, Chidambaram, Tamilnadu, India.

<sup>2</sup>Department of Computer and Information Science, Faculty of Science, Annamalai University, Chidambaram, Tamilnadu, India.

**Abstract.** The goal of this study Phishing attacks are constantly evolving, and to avoid being detected by conventional means, attackers use cutting-edge approaches. Novelty detection aims to identify previously unseen phishing attacks, including zero-day threats and sophisticated evasion tactics. Phishing attacks continue to pose significant threats to cybersecurity, exploiting human vulnerabilities and developing quickly to avoid being detected by conventional methods. In response to these challenges, this literature survey presents a comprehensive review of phishing website detection techniques, focusing on novel approaches and the latest advancements in the field. It explores dynamic analysis, real-time monitoring, and anomaly detection techniques to keep pace with the ever-changing phishing landscape. The survey addresses the persistent issue of imbalanced datasets by presenting effective strategies for handling data from significantly more legitimate websites than phishing sites. It advocates for data augmentation, cost-sensitive learning, and domain adaptation to improve the accuracy and generalization of detection models. By highlighting the latest advancements and addressing key challenges, the review contributes to building robust and resilient phishing detection frameworks that safeguard users and organizations in the constantly evolving cyber threat landscape.

**Keywords:** Phishing, Legitimate, Websites detection, Anomaly Detection, Imbalanced Datasets.

## 1 Introduction

Attacks using phishing techniques have become one of the most pervasive and sneaky cyber hazards, affecting people, businesses, and vital infrastructure all over the world. These attacks exploit social engineering techniques to deceive users into disclosing sensitive information, compromising their digital security and privacy. With phishing tactics evolving continuously, the development of effective and adaptive detection of

---

\* Corresponding author: [subaphdscholar@gmail.com](mailto:subaphdscholar@gmail.com)

website phishing techniques has become imperative to safeguard users from falling victim to these malicious schemes. In response to the escalating sophistication of phishing attacks, this literature survey presents a comprehensive examination of cutting-edge Phishing Website Detection Techniques, with a particular focus on incorporating novelty detection 1. The primary objective of this survey is to explore the latest advancements in the field and identify novel approaches that address the limitations of existing detection methods.

Traditional phishing detection mechanisms often rely on static rule-based systems or signature-based algorithms, which can struggle to keep pace with the dynamic nature of phishing attacks. Zero-day attacks and novel evasion techniques are frequently deployed by attackers to circumvent these conventional systems, demanding a paradigm shift towards more adaptive and innovative detection strategies. The central theme of this literature survey is the identification of previously unseen phishing attacks – a critical challenge in modern cybersecurity. It delves into dynamic analysis techniques that scrutinize websites in real-time, enabling the detection of zero-day attacks and rapidly evolving phishing campaigns. Furthermore, it explores anomaly detection methodologies that identify unusual patterns and behaviors, providing a proactive defense against novel phishing attempts.

One of the persistent hurdles in phishing website detection is the issue of imbalanced datasets. The overwhelming majority of legitimate websites can overshadow the limited number of phishing sites, leading to biased detection models that underperform when identifying phishing threats. To address this concern, researcher investigates data balancing techniques, cost-sensitive learning, and domain adaptation methods; all aimed at improving detection accuracy and generalization across many phishing attacks types. Recognizing the human factor in phishing attacks, this literature survey also delves into user-centric approaches. User feedback, behavior analysis, and crowd sourced intelligence play pivotal roles in developing detection systems that align with users' perspectives and augment overall security measures 2.

As the adoption of encrypted communication protocols, such as HTTPS, increases, attackers leverage encryption to conceal malicious activities. To counter this trend, it explore techniques for inspecting encrypted traffic and analyzing encrypted content, enabling the detection of phishing websites even within encrypted communications. Timeliness is of paramount importance in the detection of rapidly emerging phishing threats. Researchers emphasize the significance of real-time analysis and propose the integration of dynamic data sources, ensuring that detection models remain up-to-date and responsive to evolving attack vectors 3.

Finally, fostering collaboration within the security community is crucial in the battle against phishing attacks. Knowledge sharing, collective defense, and information exchange among researchers, organizations, and cybersecurity professionals play a pivotal role in staying ahead of the ever-evolving threat landscape. By undertaking this literature survey, our aim to present a thorough summary of the state-of-the-art Phishing Website Detection Techniques, emphasizing novelty detection and adaptive approaches. This survey will contribute to the advancement of cybersecurity by equipping practitioners and researchers with the knowledge and tools to detect and thwart sophisticated phishing attacks effectively. Ultimately, our collective efforts in developing robust and dynamic phishing detection mechanisms will enhance the security posture of users and organizations, safeguarding sensitive information from falling into the hands of malicious actors.

## 2 Background Study

### 2.1 Phishing Techniques

Phishing techniques are deceptive strategies employed by cybercriminals to trick individuals into revealing sensitive information, such as login credentials, financial details, or personal data. These techniques exploit human psychology, social engineering, and technical vulnerabilities to create convincing and trustworthy-looking traps. Here are some common phishing techniques used by attackers:

**Deceptive Emails.** Attackers send phony emails that seem to be from reliable sources, such as banks, online services, or government agencies. These emails often use urgency, fear, or enticing offers to prompt recipients to click on malicious links or download infected attachments 4.

**Spoofed Websites.** Phishers create fake websites that closely resemble legitimate ones, using similar domain names, logos, and designs. Victims are lured to these sites, where they may unknowingly enter their credentials or provide personal information, which the attackers capture 5.

**Social Engineering.** Social engineering techniques are frequently utilized in phishing attacks to manipulate victims into divulging sensitive information. The attackers could pose as someone the victim knows, such as a friend, coworker, or family member, to build trust and credibility 6.

**Pretexting.** Attackers create a fabricated scenario or pretext to trick individuals into revealing information. For example, they might impersonate IT support and claim that the victim's account has been compromised, prompting the victim to provide their login credentials 7.

**URL Manipulation.** Phishers may use URL manipulation techniques to hide malicious URLs within seemingly harmless ones. For instance, they might use URL shorteners or misspelled domain names to redirect victims to phishing websites 8.

**Credential Harvesting.** Phishing attempts seek to collect from victims confidential data such as login passwords. Once attackers obtain this information, Attackers are capable of obtaining illegal access to several accounts and systems 9.

**Email Spoofing.** Phishers use email spoofing techniques to alter the sender's address, making the message appear as if it comes from a legitimate source. This manipulation aims to deceive recipients into trusting the authenticity of the email 10.

**Voice and SMS Phishing.** Phishers may use vishing (voice phishing) or smishing (SMS phishing) techniques to deceive victims over phone calls or text messages, respectively, tricking them into divulging private details 11.

As phishing techniques continue to evolve, individuals and organizations need to stay informed about the latest threats and implement robust cybersecurity measures to protect against these deceptive attacks. User education, strong authentication mechanisms, and advanced email filtering are crucial in mitigating the risks posed by phishing attempts.

### 2.2 Phishing Detection Techniques

Phishing detection techniques aim to identify and identifying trustworthy websites from bogus ones may help you avoid phishing attempts. These techniques use various methods, including machine learning, data analysis, behavioral analysis, and website reputation assessment. Here are some common phishing detection techniques:

**Machine Learning Algorithms.** Machine learning models, such as decision trees, random forests, support vector machines (SVM), and deep neural networks, can be trained on large datasets of known phishing and legitimate websites. These models learn patterns and features indicative of phishing websites and can make accurate predictions on unseen instances 12.

**Website Content Analysis.** Phishing websites often contain specific characteristics that differentiate them from legitimate sites. Content analysis techniques examine website content, including HTML tags, URL structures, and text content, to identify suspicious elements that may indicate phishing 13.

**URL Analysis.** URL-based detection techniques inspect web addresses to identify irregularities, such as misspellings, subdomain anomalies, or the presence of foreign characters, which are common in phishing URLs 14.

**Blacklists and Whitelists.** Maintaining lists of known phishing websites (blacklists) and trusted legitimate sites (whitelists) is a straightforward approach to detecting phishing. If a website is found on the blacklist, it is blocked or flagged as suspicious 15.

**Email Authentication.** Phishing often starts with deceptive emails. Email authentication techniques, such as SPF (Sender Policy Framework) and DKIM (DomainKeys Identified Mail), verify the authenticity of the sender's domain and help identify spoofed emails 16.

**Website Certificate Verification.** SSL certificates are used to secure website connections. Phishing websites may use invalid or self-signed certificates. Verifying the authenticity of SSL certificates helps detect potential phishing attempts 17.

**Natural Language Processing (NLP).** NLP techniques analyze the language used in emails, URLs, and website content to identify phishing indicators, including suspicious grammar, vocabulary, or context 18.

**Website Reputation Services.** Utilizing reputation services or databases that track the historical behavior of websites can help identify newly registered or previously flagged phishing domains 19.

**Heuristics and Rules.** Phishing detection systems can be equipped with predefined heuristics and rules that look for specific patterns or characteristics commonly associated with phishing attacks.

**Collaborative Phishing Intelligence.** Sharing threat intelligence and collaborating with other organizations and security communities can improve the overall detection and prevention of phishing attacks.

To enhance the effectiveness of phishing detection, a combination of these techniques is often used in a layered defense strategy. By continually updating and refining detection methods, Security experts can prevent people and businesses from falling prey to these false assaults by staying ahead of the ever emerging phishing threats.

### 3 Literature Review

Several researchers have contributed significantly to the field of phishing website detection, developing various techniques and methodologies to combat the ever-evolving threat landscape. In this section, let's examine some of the well-known studies that investigated phishing detection techniques, focusing on their strengths and limitations.

**Jain Ankit Kumar et al. (2022)** 20: Jain and colleagues presented a machine learning-based approach to detect phishing websites using a combination of URL analysis, content analysis, and website structural features. Their study achieved promising results in differentiating between legitimate and phishing websites. However, the model's

performance was hindered by the lack of real-time analysis and the inability to identify novel phishing attacks.

**Jalil et al. (2022)** 21: Jalil et al. introduced an ensemble learning technique combining multiple machine learning classifiers to improve the robustness of phishing website detection. The approach showed promise in handling many phishing attacks types but did not incorporate dynamic analysis, making it vulnerable to emerging threats.

**Ramana et al. (2021)** 22: Ramana and team proposed a user-centric phishing detection approach that incorporated user behavior analysis and user feedback to enhance detection accuracy. Their study demonstrated that involving users in the detection process improved the overall effectiveness of the system. However, the model's generalization to novel phishing attacks remained a challenge.

**Harinahalli et al. (2021)** 23: Harinahalli et al. presented a real-time phishing detection framework based on deep learning techniques and dynamic analysis. Their system utilized visual similarity and content rendering to identify zero-day phishing attacks effectively. The research demonstrated promising results in handling novel threats, but there were limitations in terms of scalability and resource consumption.

**Tang et al. (2021)** 24: Sahingoz and colleagues proposed a hybrid approach that combined rule-based systems with machine learning classifiers to detect phishing websites. Their study achieved high detection rates for traditional phishing attacks, but the model's reliance on static rules limited its ability to adapt to emerging tactics.

In comparison to the existing literature, our proposed literature survey aims to bridge the gaps and limitations observed in prior works. By focusing on novelty detection and exploring dynamic analysis, user-centric approaches, and timely updates, our review seeks to offer a comprehensive understanding of the latest advancements in detecting phishing website techniques. Endeavour to provide valuable insights to researchers, practitioners, and security professionals to strengthen cybersecurity measures against the constantly evolving threat of phishing attacks.

## 4 Materials and Methods

### 4.1 Dataset

Phishing website detection often relies on large datasets containing examples of both legitimate and phishing websites. These datasets are used to train machine learning models and evaluate the performance of detection algorithms. Some popular datasets used in phishing website detection research include in [Table 1](#).

### 4.2 Phishing Detection Techniques Analysis

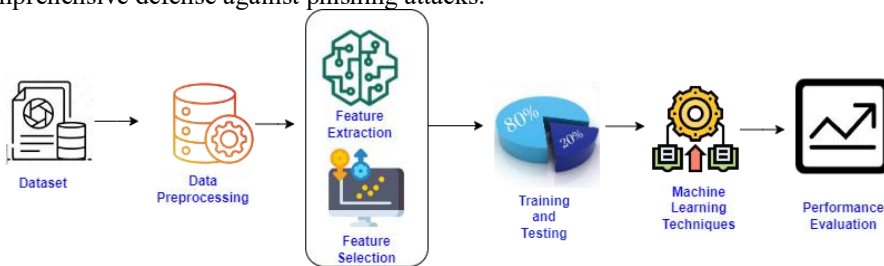
#### 4.2.1 Machine Learning-based Detection Technique

These techniques are widely used in phishing detection due to their ability to identify complex patterns and features indicative of phishing websites. These techniques leverage historical data to train models that can distinguish between legitimate and phishing websites  
**Error! Reference source not found..**

**Table 1.** Phishing and legitimate website datasets

S.No.	Dataset Name and Sources	Size of Dataset	Type	Remarks
1	Alexa 25 ( <a href="https://www.alexacom/">https://www.alexacom/</a> )	1 Million URLs	Website	Legitimate URLs
2	PhishTank 26 ( <a href="https://phishtank.org/">https://phishtank.org/</a> )	68,40,198 URLs	Website	Valid Phishing URLs
3	OpenPhish 27 ( <a href="https://openphish.com/">https://openphish.com/</a> )	11,100,075 URLs	Website	Valid Phishing URLs
4	University of california irvine ( <a href="http://www.uci.edu/">http://www.uci.edu/</a> )	2,14,748 URLs	Published Dataset	Valid Phishing URLs
5	Mendeley 28 ( <a href="https://data.mendeley.com/">https://data.mendeley.com/</a> )	2 Millions URLs	Published Dataset	Phishing + Legitimate URLs
6	CommonCrawl <b>Error!</b> <b>Reference source not found.</b> ( <a href="https://commoncrawl.org/">https://commoncrawl.org/</a> )	940 Millions URLs	Website	Legitimate URLs
7	Kaggle <b>Error!</b> <b>Reference source not found.</b> ( <a href="http://www.kaggle.com/">www.kaggle.com/</a> )	11,430 URLs	Published Dataset	Phishing + Legitimate URLs
8	Majestic Million <b>Error!</b> <b>Reference source not found.</b> ( <a href="https://majestic.com/">https://majestic.com/</a> )	Million URLs	Website	Legitimate URLs
9	Phishstorm ( <a href="https://research.aalto.fi/en/datasets/phishstorm-phishing-legitimate-url-dataset">https://research.aalto.fi/en/datasets/phishstorm-phishing-legitimate-url-dataset</a> )	96,018 URLs	Published Dataset	Phishing + Legitimate URLs
10	DMOZ ( <a href="https://dmz-odp.org/">https://dmz-odp.org/</a> )	3,861,202 URLs	Website	Legitimate URLs

Fig.1, Shows the machine learning-based phishing detection systems require high-quality training data, on-going updates to stay current with emerging threats, and careful consideration of potential biases in the data. Additionally, the deployment of machine learning models should be accompanied by other security measures to create a comprehensive defense against phishing attacks.



**Fig.1.** Machine Learning-based Detection Technique

#### 4.2.2 Deep Learning based Detection Technique

Deep learning-based detection techniques in phishing leverage neural networks with multiple layers to learn complex patterns and features from raw data in Fig.2. These techniques have shown promise in various aspects of phishing detection, including URL

analysis, content analysis, and behavior-based detection **Error! Reference source not found.** Here are some common deep learning-based detection methods in phishing:

**Convolutional Neural Networks (CNNs).** CNNs are often utilized for image recognition applications, but they can also be applied to analyze website content and URL structures. In phishing detection, CNNs can extract features from website screenshots, HTML content, or URL strings to identify patterns indicative of phishing **Error! Reference source not found.**

**Recurrent Neural Networks (RNNs).** RNNs are suitable for sequential data, making them useful for analyzing website behavior or email sequences in phishing attacks. They can learn temporal dependencies and identify unusual user interactions or content patterns that may indicate phishing attempts **Error! Reference source not found.**

**Long Short-Term Memory (LSTM) Networks.** LSTMs are a type of RNN that can better capture long-term dependencies in sequential data. They are commonly used in phishing detection to analyze user behavior, such as login sequences, form submissions, or mouse movements **Error! Reference source not found.**

**Deep Learning for URL Analysis.** Deep learning models can be used to analyze the structure and components of URLs. Techniques like character-level or word-level embedding's, along with deep neural networks, can help detect anomalies in URL strings and distinguish phishing URLs from legitimate ones **Error! Reference source not found.**

**Natural Language Processing (NLP) in Content Analysis.** Phishing emails and website content often contain suspicious language patterns. Deep learning models, combined with NLP techniques, can identify phishing indicators, such as misspellings, grammar mistakes, or unusual phrasing **Error! Reference source not found.**

**Transfer Learning for Phishing Detection.** Transfer learning allows pre-trained models from related tasks, such as image recognition or language processing, to be fine-tuned for phishing detection. This strategy can be useful when labelled phishing data is limited **Error! Reference source not found.**

**Ensemble Deep Learning Models.** Ensemble techniques, such as stacking or blending multiple deep learning models, can improve detection accuracy and robustness by combining the strengths of different architectures **Error! Reference source not found.**

**Online Learning for Real-time Detection.** Deep learning models can be adapted for online learning, enabling real-time updates and continuous improvement to respond rapidly to emerging phishing threats **Error! Reference source not found.**

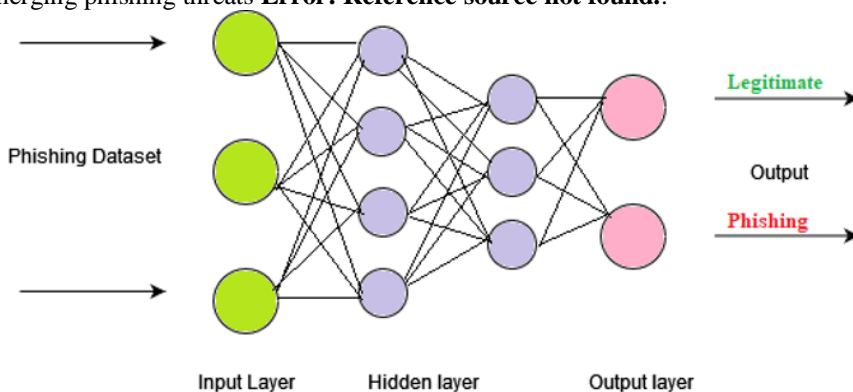


Fig.2. Deep Learning-based Detection Technique

### 4.3 Anomaly detection methods

Anomaly detection methods play a crucial role in phishing website detection by identifying deviations from normal patterns or behaviors. Anomalous activities in web content, user behavior, or network traffic may indicate the presence of phishing attempts. These methods aim to detect novel and sophisticated phishing attacks that might not be recognized by traditional rule-based or signature-based systems. Machine learning techniques like anomaly detection look for instances or data points that drastically differ from the dataset as a whole. The assumption is that anomalies represent unusual or suspicious behavior that requires further investigation. In the context of phishing website detection, anomalies may represent malicious websites that exhibit characteristics different from legitimate websites.

Anomaly detection methods in phishing website detection require the extraction of relevant features from the web content, URL, or user interactions. These features serve as input data to the anomaly detection algorithms. Common features might include URL components, website structure, textual content, image metadata, and user interaction patterns. Overall, anomaly detection methods provide a valuable approach to detecting novel and emerging phishing attacks by identifying patterns that differ from normal behavior. When combined with other detection techniques, such as rule-based and machine learning-based methods, anomaly detection contributes to more comprehensive and effective phishing website detection systems.

### 4.4 Addressing Imbalanced Datasets

Addressing imbalanced datasets is an important aspect of phishing website detection, as real-world datasets often have a significantly the quantity of trustworthy websites compared to the relatively smaller number of fraudulent website. Imbalanced data poses challenges for machine learning algorithms, as they may become biased towards the majority class (legitimate websites) and perform poorly in detecting the minority class (phishing websites).

**Data Resampling.** Data resampling techniques are commonly used to balance the dataset by either increasing the minority class or decreasing the number of majority class. Two primary approaches are:

**Oversampling.** Generating synthetic instances of the underclass by duplicating existing samples or using techniques like SMOTE (Synthetic Minority Over-sampling Technique) to create synthetic examples that are similar to existing phishing instances.

**Under sampling.** Reducing occurrences in most of the class as a percentage to match the number of the minority class, either randomly or using more sophisticated techniques like Tomek links or Cluster Centroids.

### 4.5 Ensemble Techniques

Ensemble methods, such as Bagging and Boosting, can help improve the performance of imbalanced datasets. Bagging combines multiple classifiers trained on different subsets of the data, reducing the risk of overfitting to the majority class. Boosting, on the other hand, increases the weights of occurrences that are incorrectly categorised, emphasizing the minority class more during the training phase. Addressing imbalanced datasets is critical to ensure that phishing website detection systems effectively identify both legitimate and phishing websites. By applying these techniques, researchers and practitioners can improve



the performance and accuracy of detection models, making them more robust against the difficulties caused by unbalanced data.

#### 4.6 Novelty Detection

Novelty detection in phishing website detection refers to the ability of a system to identify and handle previously unseen or novel phishing attacks. Phishing attacks are continuously evolving, and attackers use various tactics to avoid conventional detecting techniques. Novelty detection techniques aim to overcome this limitation by focusing on detecting unknown or zero-day phishing threats, which have not been encountered before.

**Detecting Zero-Day Attacks.** Zero-day attacks refer to newly emerging phishing threats that exploit unknown vulnerabilities or weaknesses in security systems. Novelty detection methods can identify and flag these attacks even without prior knowledge of their existence.

**Dynamic Analysis and Real-Time Monitoring.** Dynamic analysis techniques, including real-time monitoring and behavior analysis, play a vital role in novelty detection. By analyzing website behavior in real-time, these methods can capture and identify novel phishing attacks as they emerge.

**Handling Encrypted Traffic.** Novelty detection methods must also address encrypted traffic since attackers can use encryption to cover up their harmful actions. Techniques like TLS/SSL interception and encrypted content analysis help identify novel phishing websites transmitted over secure channels.

**Robustness against Evasion Tactics.** Attackers may use evasion tactics to make phishing websites appear more legitimate and evade detection. Novelty detection techniques should be designed to be robust against these evasion tactics.

#### 4.7 Performance Evaluation

In our experiments, using the following metrics, researchers assess how effectively the phishing detection methods perform: True positive rate (TPR), False positive rate (FPR), precision, f-score, recall, and accuracy (ACC). According to the equations below, the metrics were calculated.

*False negative rate (FNR).* The number of phishing websites that are improperly categorised is shown below in a formula (1).

$$FNR = \frac{FNR}{FNR+TPR} \quad (1)$$

*False positive rate (FPR).* In this formula below, FPR stands for the fraction of legal websites that are mistakenly labelled as phishing sites. FP called as phishing website's. Where TN as legitimate websites identified accurately as shown in formula (2).

$$FPR = \frac{FPR}{FPR+TNR} \quad (2)$$

*Recall.* The ratio of correctly predicted rumour tweets (True Positives) to all other tweets (True Positives + False Negatives) as shown in formula (3).

$$Recall. = \frac{TPR}{TPR+FNR} \quad (3)$$

*Precision.* This measures the percentage for correctly predicted rumor tweets (True Positives) to all previously identified rumour tweets (True Positives + False Positives) as shown in formula (4).

$$precision = \frac{TPR}{TPR+FPR} \quad (4)$$

*F-score.* This has precision and memory and is symmetrical. It achieved a compromise between evaluations of recall and precision as shown in formula (5).

$$f - score = 2 \times \frac{precision \times recall}{precision+recall} \quad (5)$$

*Accuracy (ACC).* ACC refers to the proportion of websites with the proper classification, including those that are legitimate websites and those that are accurately identified as phishing websites as shown in formula (6).

$$Acc. = \frac{TPR+TNR}{FPR+FNR+TPR+TNR} \quad (6)$$

## 5 Issues and Challenges

**Phishing Attack Sophistication.** Phishing attacks are becoming increasingly sophisticated, employing advanced techniques like social engineering, targeted spear-phishing, and brand impersonation. Traditional detection methods may struggle to keep up with these evolving tactics.

**Zero-Day Phishing Attacks.** Zero-day phishing attacks exploit vulnerabilities that have not been previously identified or patched. Detecting and defending against these attacks before they are widely known presents a significant challenge.

**Imbalanced Datasets.** Phishing datasets are often imbalanced, with a larger number of trustworthy websites compared to phishing websites. This imbalance can lead to biased classifiers and a higher risk of false negatives.

**Encrypted Traffic.** The growing adoption of HTTPS and encryption poses challenges in inspecting encrypted traffic for potential phishing websites, as attackers can use encryption to cover up their harmful actions.

**Evasion Techniques.** Phishers continuously develop evasion tactics to bypass detection systems. These techniques may involve obfuscation, URL shorteners, or polymorphic phishing attacks, making it challenging to identify and block phishing attempts.

**User-Centric Attacks.** Phishing attacks that target specific users or organizations, known as spear-phishing, can be difficult to detect as they are personalized and tailored to their victims.

**Cross-Platform Attacks.** Phishing attacks target various platforms, including web browsers, mobile devices, and email clients. Detecting phishing attempts consistently across different platforms is challenging due to their varying characteristics.

**Real-Time Detection.** Timely detection is crucial in mitigating the impact of phishing attacks. Real-time detection systems are needed to respond rapidly to emerging phishing threats.

**Evasion of Machine Learning-Based Detection.** Sophisticated phishers can design attacks to evade machine learning-based detection systems by crafting phishing websites that resemble legitimate sites or by generating adversarial examples.

**User Education and Awareness.** Despite sophisticated detection systems, user education and awareness remain critical. Human error, such as falling for phishing emails or social engineering tactics, can still be a significant risk factor.

## 6 Research Gap

The literature review on phishing website detection has undoubtedly shed light on the diverse array of techniques and methodologies employed in this critical field of cybersecurity. However, within this extensive body of research, there emerge certain notable research gaps that warrant further investigation and exploration.

Firstly, while the literature review extensively discusses the various detection methods, it becomes evident that a gap exists in terms of a unified approach that effectively combines multiple techniques. Many of the reviewed methods demonstrate strengths in specific scenarios, but a holistic system that harnesses the benefits of rule-based, machine learning-based, dynamic analysis, and anomaly detection approaches remains underexplored. Developing such a comprehensive framework could significantly enhance the overall accuracy and robustness of phishing website detection.

Secondly, the review highlights the consistent challenge of handling encrypted traffic, which is an emerging concern as more websites adopt secure communication protocols. However, the discussion primarily revolves around the detection of phishing within encrypted channels, rather than exploring the potential application of encryption as a defense mechanism. Investigating the feasibility of leveraging encryption techniques to secure sensitive user information from phishing attacks could be an area ripe for exploration.

Finally, the literature review provides a rich foundation of insights into phishing website detection techniques and their associated challenges. However, it underscores the importance of addressing key research gaps: developing an integrated approach, exploring encryption as a defense mechanism, refining user-centric methods, and investigating the potential of emerging technologies in real-time detection. Closing these gaps would not only contribute to the advancement of the field but also bolster the overall cybersecurity ecosystem.

## 7 Conclusion

In conclusion, the literature review on a thorough overview of the methods now in use and recent developments in the field has been given by the detection of phishing websites. The review explored various detection methods, including traditional rule-based and signature-based approaches, as well as modern machine learning-based and anomaly detection methods. The review identified several limitations in traditional detection techniques, including their susceptibility to sophisticated and novel phishing attacks. Additionally, the challenge of handling imbalanced datasets, where legitimate websites vastly outnumber phishing websites, was emphasized. One of the significant highlights of the literature review was the exploration of novelty detection techniques, which aim to identify and handle previously unseen or zero-day phishing attacks. These methods play a critical role in addressing the constantly evolving nature of phishing threats and improving the resilience of detection systems. Overall, the literature review highlighted the dynamic and evolving nature of phishing website detection, with on-going research focusing on

explainable AI, adversarial defense, privacy-preserving techniques, and cross-platform detection. In order to keep ahead of sophisticated phishing attempts and protect people and businesses from falling for these trickery tactics as the field of cybersecurity changes, ongoing research and innovation will be necessary. The insights provided in this literature review serve as a valuable resource for researchers, practitioners, and cybersecurity professionals working towards more effective and robust phishing detection systems.

## References

1. Chiew Kang Leng Kelvin Sheng Chek Yong and Choon Lin Tan 2018 A survey of phishing attacks: Their types, vectors and technical approaches. *Expert Systems with Applications*. 106: 1-20. Doi: <https://doi.org/10.1016/j.eswa.2018.03.050>
2. Rao Routhu Srinivasa and Alwyn Roshan Pais 2019 Detection of phishing websites using an efficient feature based machine learning framework. *Neural Computing and Applications*. 31(8): 3851-3873. Doi: <https://doi.org/10.1007/s00521-017-3305-0>
3. <http://www2.deloitte.com/content/dam/Deloitte/sg/Documents/risk/searisk-cyber-101-part10.pdf>
4. G J W Kathrine P M Praise A A Rose and E C Kalaivani 2019 Variants of phishing attacks and their detection techniques 3rd International Conference on Trends in Electronics and Informatics (ICOEI). 255-259. DOI: [10.1109/ICOEI.2019.8862697](https://doi.org/10.1109/ICOEI.2019.8862697)
5. Rao R S Pais A R 2019 Detection of phishing websites using an efficient feature-based machine learning framework. *Neural Comput & Applic*. 31: 3851–3873. DOI: <https://doi.org/10.1007/s00521-017-3305-0>
6. Rao R S Pais A R and Anand P 2020 A heuristic technique to detect phishing websites using TWSVM classifier. *Neural Comput & Applic* DOI: <https://doi.org/10.1007/s00521-020-05354-z>
7. S Roopak A P Vijayaraghavan and T Thomas 2019 On Effectiveness of Source Code and SSL Based Features for Phishing Website Detection. 1st International Conference on Advanced Technologies in Intelligent Control. *Environment, Computing & Communication Engineering (ICATIECE)*:172-175. DOI: [10.1109/ICATIECE45860.2019.9063824](https://doi.org/10.1109/ICATIECE45860.2019.9063824)
8. A Nakamura and F Dobashit 2019 Proactive Phishing Sites Detection. *IEEE/WIC/ACM International Conference on Web Intelligence (WI)*. 443-448 DOI: <https://doi.org/10.1145/3350546.3352565>
9. F Tajaddodianfar J W Stokes and A Gururajan 2020 Texception: A Character/Word-Level Deep Learning Model for Phishing URL Detection. *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 2857-2861 DOI: [10.1109/ICASSP40776.2020.9053670](https://doi.org/10.1109/ICASSP40776.2020.9053670)
10. K Althobaiti G Rummani and K Vaniea 2019 A Review of Human and Computer Facing URL Phishing Features. *IEEE European Symposium on Security and Privacy Workshops*. 182-191 DOI: [10.1109/EuroSPW.2019.00027](https://doi.org/10.1109/EuroSPW.2019.00027)
11. Carlo Marcelo Revoredo da Silva Eduardo Luzeiro Feitosa Vinicius Cardoso Garcia 2020 Heuristic based strategy for Phishing prediction: A survey of URL-based approach. *Computers & Security*, 101613 DOI: <https://doi.org/10.1016/j.cose.2019.101613>
12. Athulya A A and K Praveen 2020 Towards the detection of phishing attacks. *4th international conference on trends in electronics and informatics (ICOEI)(48184)*. DOI: [10.1109/ICOEI48184.2020.9142967](https://doi.org/10.1109/ICOEI48184.2020.9142967)

13. Sahar Abdelnabi Katharina Krombholz and Mario Fritz 2020 VisualPhishNet: Zero-Day Phishing Website Detection by Visual Similarity. *Association for Computing Machinery*. 1681–1698 DOI: <https://doi.org/10.1145/3372297.3417233>
14. S Haruta H Asahina and I Sasase 2017 Visual Similarity-Based Phishing Detection Scheme Using Image and CSS with Target Website Finder. *IEEE Global Communications Conference*. pp. 1-6. DOI: [10.1109/GLOCOM.2017.8254506](https://doi.org/10.1109/GLOCOM.2017.8254506)
15. M M Yadollahi F Shoeleh E Serkani A Madani and H Gharaee 2019 An Adaptive Machine Learning Based Approach for Phishing Detection Using Hybrid Features. *Web Research*. pp. 281-286 DOI: [10.1109/ICWR.2019.8765265](https://doi.org/10.1109/ICWR.2019.8765265)
16. Jain AK Gupta B B 2019 A machine learning based approach for phishing detection using hyperlinks information. *J Ambient Intell Human Comput* 10. 2015–2028 DOI: <https://doi.org/10.1007/s12652-018-0798-z>
17. J Kumar A Santhanavijayan B Janet B Rajendran and B S Bindhumadhava 2020 Phishing Website Classification and Detection Using Machine Learning. *Computer Communication and Informatics*. pp. 1-6 DOI: <https://doi.org/10.48550/arXiv.2103.12739>
18. <https://www.kdnuggets.com/2020/02/deepneural-networks.html>
19. I Saha D Sarma R J Chakma M N Alam A Sultana and S Hossain 2020 Phishing Attacks Detection using Deep Learning Approach. *Smart Systems and Inventive Technology*. pp. 1180-1185 DOI: [10.1109/ICSSIT48917.2020.9214132](https://doi.org/10.1109/ICSSIT48917.2020.9214132)
20. Jain Ankit Kumar and B B Gupta 2022 A survey of phishing attack techniques defence mechanisms and open research challenges. *Enterprise Information Systems*. 16(4): 527-565 DOI: <https://doi.org/10.1080/17517575.2021.1896786>
21. Jalil Sajjad Muhammad Usman and Alvis Fong 2022 Highly accurate phishing URL detection based on machine learning. *Journal of Ambient Intelligence and Humanized Computing*: 1-19 DOI: <https://doi.org/10.1007/s12652-022-04426-3>
22. Ramana A V K Lakshmana Rao and Routhu Srinivasa Rao 2021 Stop-Phish: an intelligent phishing detection method using feature selection ensemble. *Social Network Analysis and Mining*. 11(1): 1-9 DOI: <https://doi.org/10.1007/s13278-021-00829-w>
23. Harinahalli Lokesh Gururaj and Goutham BoreGowda 2021 Phishing website detection based on effective machine learning approach. *Journal of Cyber Security Technology*. 5(1): 1-14 DOI: <https://doi.org/10.1080/23742917.2020.1813396>
24. Tang Lizhen and Qusay H Mahmoud 2021 A survey of machine learning-based solutions for phishing website detection. *Machine Learning and Knowledge Extraction*. 3(3): 672-694 DOI: <https://hdl.handle.net/10155/1446>
25. <https://www.alexa.com/topsites>
26. <http://index.commoncrawl.org/>
27. [https://www.phishtank.com/developer\\_info.php](https://www.phishtank.com/developer_info.php)
28. <https://openphish.com/>