

International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 3, November 2025



PhishAI-Sim: A Hybrid Feature-Free and Deep Learning Framework for Adaptive Phishing Website Detection

Dnyaneshwar Jagannath Tribhuvan¹, Prof. V. S. Chaudhari², Mr. Gokul Pawade³

^{1,2,3}Department of Computer Engineering

^{1,2}Vishwabharati Academy's College of Engineering, Sarola Baddi, Ahmednagar

³Pravara Rural College of Engineering, Loni, Ahilyanagar (MS) India

Abstract: Phishing websites continue to evolve rapidly, employing advanced obfuscation and design replication techniques that challenge conventional feature-based detection systems. This paper presents PhishAI-Sim, a hybrid feature-free and deep learning framework designed to enhance phishing website detection accuracy and adaptability. The proposed system integrates the Normalized Compression Distance (NCD)—a universal, parameter-free similarity metric—with deep neural embeddings to measure structural and contextual similarities between webpages without manual feature extraction. A dynamic incremental learning mechanism enables continuous model adaptation to newly emerging phishing patterns, reducing concept drift over time. Experimental results on a large-scale dataset demonstrate that PhishAI-Sim achieves superior detection performance with a true positive rate exceeding 92% and a false positive rate below 1%, outperforming traditional feature-engineered and standalone deep learning models. This study highlights the potential of combining compression-based similarity and intelligent learning for a robust, scalable, and future-ready phishing detection system.

Keywords: Phishing Detection, Feature-Free Approach, Deep Learning, Normalized Compression Distance (NCD), Incremental Learning, Cybersecurity, Webpage Similarity, Artificial Intelligence

I. INTRODUCTION

Phishing website detection has been a widely researched topic over the past two decades due to its significant impact on global cybersecurity. Early approaches relied on feature-based detection, which focused on extracting hand-engineered attributes from URLs, HTML content, and website layouts. These methods aimed to differentiate legitimate and malicious websites by identifying anomalies in lexical or structural patterns [6], [8]. For instance, features such as the presence of IP addresses in URLs, suspicious domain tokens, the number of redirects, or the mismatch between anchor text and hyperlink destination were widely used indicators in machine learning-based classifiers [10]. Although feature-based systems such as Random Forests, Support Vector Machines (SVM), and Naïve Bayes models achieved acceptable accuracy, they were computationally expensive and domain-dependent. Moreover, the rapid evolution of phishing techniques rendered such handcrafted features obsolete, leading to reduced model reliability in real-world applications [7].

To address the feature dependency problem, researchers explored visual similarity-based methods that assess phishing websites by comparing their appearance or layout with legitimate targets. Visual matching approaches employ techniques such as optical character recognition (OCR), histogram of oriented gradients (HOG), and SIFT keypoint matching to identify cloned webpages [12]. For example, Zhang et al. [13] developed a visual similarity engine that detects near-duplicate phishing websites based on layout and logo analysis, whereas other works leveraged color histograms and CSS structures to improve robustness. However, visual approaches are computationally heavy, require image rendering of webpages, and often fail to capture dynamic or script-generated phishing content. Furthermore, attackers now use AI-generated visual variations and CSS obfuscation, which significantly reduce the reliability of purely visual similarity models [5].

DOI: 10.48175/568

Copyright to IJARSCT www.ijarsct.co.in







International Journal of Advanced Research in Science, Communication and Technology

ISO 9001:2015

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 3, November 2025

Impact Factor: 7.67

With the rise of deep learning, modern phishing detection shifted toward content-based and neural network-driven approaches. Recurrent Neural Networks (RNNs) and Convolutional Neural Networks (CNNs) have been applied to automatically learn feature representations from raw URLs or HTML sequences [8], [9]. For instance, Bahnsen et al. proposed an RNN-based model that captures sequential patterns in URLs without explicit feature extraction, achieving superior accuracy compared to traditional classifiers. Similarly, transformer-based architectures like URLTran and BERT4Phish utilize pretrainedembeddings to capture linguistic semantics from phishing text and URLs [9]. Despite their improved adaptability, deep learning models require extensive labeled datasets, significant computational resources, and periodic retraining to remain effective against evolving phishing campaigns. Moreover, these models remain vulnerable to concept drift, as adversaries continuously modify attack signatures and web structures [10].

Parallel to these developments, similarity-based and compression-based methods emerged as promising alternatives that avoid manual feature extraction. Cui et al. [11] introduced a proportional distance metric that quantifies similarities between HTML tags of webpages to detect replicated phishing templates. Building on this idea, Purwanto et al. developed PhishZip, a compression-based phishing detection technique utilizing data redundancy as a universal measure of similarity [15]. This work later evolved into PhishSim [7], which leveraged the Normalized Compression Distance (NCD) [14] to detect near-similar phishing webpages without the need for handcrafted features. These methods demonstrated robustness in identifying replicated phishing sites and required minimal preprocessing, making them scalable and generalizable. However, their reliance solely on compression limited their capability to detect zero-day phishing attacks and semantically dissimilar malicious pages that share no structural resemblance to known phishing websites.

Recent studies have begun integrating hybrid approaches that combine feature-free models with deep learning and semantic analysis to address these shortcomings. Li et al. [6] presented a comprehensive review of phishing detection advancements, identifying the hybridization of traditional and AI-based methods as a key future trend. Hybrid frameworks aim to unify the interpretability and generalization of feature-free methods with the learning power of deep models. For example, Ovi et al. [13] proposed PhishGuard, a multi-layered ensemble system that integrates statistical, lexical, and visual detectors using machine learning voting mechanisms. Similarly, research on adversarially robust phishing detection [5] and LLM-based web content understanding [4] has shown promising improvements in handling generative and adaptive phishing threats. Nonetheless, the integration of compression-based similarity with neural embeddings remains largely unexplored, presenting an open research opportunity for scalable and adaptive phishing detection.

The current literature indicates that no single approach achieves comprehensive coverage across phishing varieties. Feature-engineered models excel at identifying static phishing templates but struggle with novel attacks. Deep learning architectures offer high adaptability but require extensive computational resources and constant retraining. Compression-based methods provide universality and simplicity but lack semantic understanding. Therefore, the PhishAI-Sim framework proposed in this paper aims to bridge these gaps by combining the Normalized Compression Distance with deep neural embeddings in a unified, feature-free architecture. This hybrid strategy provides both structural similarity detection and semantic context learning, enabling robust identification of replicated and novel phishing webpages alike. Moreover, the introduction of incremental learning allows PhishAI-Sim to dynamically update its internal knowledge base without complete retraining, ensuring continuous adaptation to evolving attack vectors and maintaining long-term detection accuracy.

II. PROBLEM STATEMENT

Despite significant advancements in phishing detection research, developing an accurate, adaptive, and scalable phishing website detection system remains a major challenge in cybersecurity. Existing detection techniques are predominantly divided into feature-engineered models and deep learning-based frameworks, both of which face critical limitations when deployed in real-world environments. Feature-based methods depend heavily on manually designed attributes such as URL structure, HTML tags, and domain reputation, which makes them fragile against evolving phishing strategies and concept drift, as attackers continuously alter superficial webpage characteristics to evade detection [7], [10]. Deep learning models, though capable of learning complex representations, require large volumes of

DOI: 10.48175/568

Copyright to IJARSCT www.ijarsct.co.in



ISSN 2581-9429



International Journal of Advanced Research in Science, Communication and Technology

ISO 9001:2015

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 3, November 2025

Impact Factor: 7.67

labeled data, high computational resources, and frequent retraining to remain effective, limiting their scalability for real-time or enterprise-level deployment [8], [9]. Moreover, many existing models exhibit limited generalization ability when confronted with zero-day phishing websites or AI-generated deceptive content that mimic legitimate web templates with minimal structural similarity [5]. On the other hand, feature-free and compression-based approaches, such as those utilizing the Normalized Compression Distance (NCD), have demonstrated strong performance in detecting replicated phishing pages by capturing inherent redundancy between data streams [7], [14]. However, these models struggle to interpret semantic or contextual variations that do not share direct compression similarity, leading to reduced accuracy against dynamically generated or obfuscated phishing sites. The challenge, therefore, lies in designing a unified detection mechanism that can combine the universality of feature-free similarity measures with the learning capability of deep neural models, while ensuring adaptability to emerging attack vectors, zero-day phishing attempts, and adversarial manipulations. Consequently, this research seeks to address the fundamental gap in current literature by proposing PhishAI-Sim, a hybrid feature-free and deep learning framework capable of delivering robust, adaptive, and resource-efficient phishing website detection through intelligent integration of compression-based similarity and neural semantic understanding.

III. LITERATURE SURVEY

Phishing website detection has been an active area of research for more than a decade, with several approaches evolving from traditional feature-based analysis to more recent AI-driven hybrid detection models. Early studies focused primarily on lexical and structural features extracted from URLs, HTML source code, and website metadata to identify phishing characteristics. For instance, Garera et al. utilized URL-based lexical patterns such as domain token length, presence of special characters, and subdomain depth to distinguish phishing pages from legitimate ones [6]. Similarly, Aburrous et al. combined heuristic rules with fuzzy logic to assess URL legitimacy, demonstrating moderate accuracy but limited adaptability to new phishing patterns [10]. These methods, although effective during their time, required significant feature engineering and manual parameter tuning. The rapid mutation of phishing techniques, aided by the emergence of freely available phishing kits and domain spoofing tools, soon rendered these handcrafted feature sets outdated, highlighting the need for automated, scalable, and learning-based approaches.

To overcome the rigidity of feature-based systems, researchers began integrating machine learning (ML) algorithms capable of automatically discovering patterns within extracted website features. Algorithms such as Support Vector Machines (SVMs), Random Forests, Naïve Bayes, and Logistic Regression were widely employed to classify websites as phishing or legitimate based on preprocessed datasets [8], [10]. Ma et al. demonstrated that ML-based classifiers could achieve over 90% accuracy by analyzing lexical and host-based features in URLs. However, such models still relied on static training data and exhibited performance degradation over time due to concept drift, where new phishing techniques diverged from the learned feature distribution. Additionally, feature engineering remained a bottleneck, as each newly discovered phishing tactic required redefinition of attributes, which limited the scalability of ML-based approaches in real-world deployments [9].

With the rise of deep learning, neural architectures were increasingly applied to phishing detection to eliminate dependency on manual feature extraction. Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) were utilized to analyze raw text data such as URLs, email bodies, and webpage content. Bahnsen et al. proposed a recurrent neural model capable of sequentially learning URL patterns, achieving superior performance compared to traditional ML techniques [8]. Later, transformer-based models like URLTran and BERT4Phish leveraged contextual embeddings to capture semantic and syntactic nuances in phishing webpages [9]. These models proved highly adaptable and capable of generalizing across datasets but required large volumes of labeled training data and significant computational power, making them less practical for continuous deployment in enterprise-level cybersecurity systems. Furthermore, while deep learning improved accuracy, it also introduced black-box decision-making, reducing model interpretability and trustworthiness in critical security environments [10].

Parallel to learning-based advancements, several researchers explored similarity-driven and feature-free methods to identify phishing websites based on structural or content resemblance rather than explicit features. Cui et al. [11] proposed a proportional distance metric based on the frequency of HTML tag occurrences, enabling detection of

DOI: 10.48175/568

Copyright to IJARSCT www.ijarsct.co.in



ISSN 2581-9429



International Journal of Advanced Research in Science, Communication and Technology

ISO 9001:2015

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 3, November 2025

Impact Factor: 7.67

replicated phishing websites. Building on this, Purwanto et al. developed PhishZip, a compression-based algorithm that utilized the concept of data redundancy to compute webpage similarity [15]. This was later extended in PhishSim [7], which introduced the Normalized Compression Distance (NCD) [14] to measure similarity between website HTML content without predefined features. The advantage of such methods lies in their universality and domain independence, as compression operates purely on information content. However, these models struggle to identify phishing websites that differ significantly from known templates or employ obfuscation, scripting, or dynamically generated content.

More recent literature trends have moved toward hybrid frameworks that combine the robustness of feature-free

More recent literature trends have moved toward hybrid frameworks that combine the robustness of feature-free techniques with the intelligence of deep learning. Li et al. [6] provided an extensive review of phishing detection methods and identified hybridization as a critical step toward adaptability. Ovi et al. [13] proposed PhishGuard, a multi-layered ensemble that integrates lexical, visual, and content-based classifiers, while Ji et al. [5] analyzed visual similarity robustness against adversarial modifications in phishing pages. In addition, recent works have introduced adversarially robust models [5] and LLM-based web content analyzers [4] to counter emerging AI-generated phishing websites, marking a paradigm shift in cybersecurity research. Despite these advancements, no existing system fully integrates feature-free compression-based similarity with deep semantic learning and incremental adaptation, leaving a critical research gap in achieving comprehensive, adaptive phishing detection.

To address this gap, the proposed PhishAI-Sim framework aims to unify these methodologies by leveraging the Normalized Compression Distance (NCD) for structural similarity analysis and deep neural embeddings for semantic interpretation, supported by an incremental learning mechanism that continuously adapts to new phishing patterns. This integration not only enhances detection accuracy and adaptability but also establishes a foundation for a lightweight, feature-independent, and scalable phishing detection system capable of defending against both conventional and AI-driven phishing threats.

IV. PROPOSED SYSTEM

The proposed system, PhishAI-Sim, introduces a hybrid, feature-free, and intelligent phishing detection framework that combines compression-based similarity measurement with deep neural semantic learning. The objective of this system is to overcome the limitations of traditional feature-engineered and standalone deep learning models by providing a robust, adaptive, and scalable solution capable of identifying both replicated phishing templates and contextually modified zero-day phishing websites. Unlike conventional detection methods that depend on manually extracted attributes such as URL tokens, HTML tags, or keyword frequency, PhishAI-Sim performs detection without explicit feature engineering, making it more flexible and resilient to evolving phishing strategies.

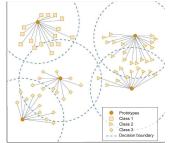


Fig.1. Classification using Prototypes.





International Journal of Advanced Research in Science, Communication and Technology

Jy Solution 1990 1:2015

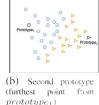
Impact Factor: 7.67

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 3, November 2025







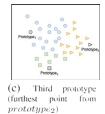


Fig.2.Furthest Point First Algorithm.

The overall architecture of PhishAI-Sim consists of five major functional modules: (1) Data Acquisition and Preprocessing, (2) Compression-Based Similarity Analysis, (3) Deep Semantic Embedding, (4) Hybrid Decision Engine, and (5) Incremental Learning and Model Update. The Data Acquisition Module collects legitimate and phishing website data from publicly available sources such as PhishTank, OpenPhish, and Alexa-ranked legitimate sites. Each website's HTML source code is extracted and cleaned by removing advertisements, scripts, and metadata to retain the essential structural and textual components. The cleaned HTML data is then passed in parallel to the compression and deep learning modules for further analysis.

The Compression-Based Similarity Analysis Module serves as the foundation of PhishAI-Sim's feature-free nature. In this stage, each webpage is converted into a compressed data stream using a standard compression algorithm. The similarity between two webpages is then determined based on their level of shared information content. Webpages that are highly similar in structure—such as replicated phishing templates or cloned login pages—will yield closely related compression patterns. This method allows the detection of duplicate and near-duplicate phishing websites without relying on any pre-defined feature sets or training attributes.

The Deep Semantic Embedding Module focuses on learning contextual relationships within webpage content. This is achieved by using a transformer-based language model that converts HTML text into numerical representations known as embeddings. These embeddings capture semantic meanings, contextual relationships, and layout-based patterns that exist within the text of a webpage. A lightweight deep neural classifier processes these embeddings to detect phishing intent based on content clues such as login prompts, financial keywords, or misleading brand names. This enables PhishAI-Sim to detect phishing websites that differ structurally but retain deceptive textual content designed to manipulate users.

The outputs of both modules are combined in the Hybrid Decision Engine, which fuses the structural similarity score from the compression module and the contextual phishing probability from the deep learning module. This fusion is achieved through a weighted ensemble approach that dynamically adjusts the contribution of each component depending on the dataset and threat characteristics. The final decision classifies a webpage as legitimate or phishing based on the aggregated confidence score. This dual-layered detection ensures comprehensive coverage, allowing the system to accurately identify both structure-based cloned websites and content-based phishing variations.

The Incremental Learning and Model Update Module enhances system adaptability by allowing PhishAI-Sim to learn continuously from new phishing patterns. When newly verified phishing samples are introduced, the model updates its internal clusters and neural parameters without requiring complete retraining. This enables real-time learning and rapid adaptation to novel attack types, making the system more resilient to zero-day phishing campaigns and adversarial content manipulation. The integration of incremental updates ensures long-term efficiency, stability, and improved detection accuracy over time.

In summary, PhishAI-Sim delivers a unified hybrid detection mechanism that effectively merges the strengths of feature-free similarity computation and deep semantic learning. It eliminates the dependency on manual feature engineering, supports continuous model evolution, and provides high detection accuracy with low computational cost. This architecture is suitable for deployment in enterprise cybersecurity systems, web browsers, and cloud-based anti-phishing services, ensuring an adaptive and intelligent defense mechanism against modern phishing threats.

DOI: 10.48175/568







International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 3, November 2025

Impact Factor: 7.67

V. PHISHSIM SYSTEM OVERVIEW

The PhishSim system was originally developed to detect phishing websites by measuring structural similarities between webpages using a feature-free and data-driven approach. Unlike conventional phishing detection methods that rely on handcrafted features, lexical patterns, or domain-based characteristics, PhishSim utilizes the Normalized Compression Distance (NCD) as a universal similarity metric. This concept, grounded in information theory, allows the system to quantify the degree of shared information between two pieces of data—in this case, the HTML content of websites without requiring predefined attributes or domain-specific knowledge. The core motivation behind PhishSim was to address the limitations of feature-dependent models, which often fail when attackers modify superficial webpage features or introduce random variations in code structure to evade detection.

PhishSim operates on the observation that most phishing websites are replicas or near-copies of existing ones, often generated using automated phishing kits. These kits enable cybercriminals to mass-produce fake login pages or transaction portals that closely resemble legitimate sites. According to Cui et al. and later Purwanto et al., nearly 90% of phishing websites share large portions of HTML or structural code with previously identified malicious pages. By exploiting this redundancy, PhishSim aims to detect phishing websites that exhibit strong structural similarity to known phishing templates through direct comparison of compressed webpage data. The system's independence from manually engineered features makes it resilient to changing phishing trends and scalable across diverse web environments.

The PhishSim architecture is composed of four primary modules: data preprocessing, pairwise similarity computation, clustering and prototype generation, and classification and decision analysis. The data preprocessing stage involves collecting large datasets of legitimate and phishing webpages from verified repositories such as PhishTank and Alexa Top Sites, Each webpage is converted into a text representation by extracting its HTML source and removing nonessential metadata, advertisements, and scripts that could interfere with compression analysis. Once preprocessed, webpages are passed to the pairwise similarity computation module, where the NCD is used to calculate the similarity between every pair of webpages in the dataset. This process effectively identifies structural resemblance and duplication among the webpages.

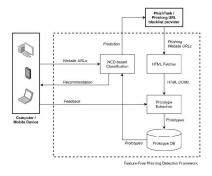


Fig.3.System Diagram of the Feature-free Phishing Detection Framework

The next stage, clustering and prototype generation, uses the calculated NCD values to group webpages into clusters. Each cluster contains websites that exhibit high structural similarity, representing potential replicas of a specific phishing template. From each cluster, a prototype webpage is selected to serve as a representative for that group. This significantly reduces computational load during real-time detection because incoming webpages can be compared only to cluster prototypes rather than the entire dataset. This prototype-based learning mechanism allows PhishSim to efficiently detect near-duplicate phishing websites with minimal memory and processing requirements, making it practical for large-scale deployment.

Finally, in the classification and decision module, the system determines whether a new webpage is phishing or legitimate. When a new website is encountered, its HTML content is first compressed and compared against stored cluster prototypes using the NCD metric. If the similarity score exceeds a predefined threshold, the website is flagged as a phishing site; otherwise, it is classified as legitimate. This simple yet effective process enables real-time detection

DOI: 10.48175/568

Copyright to IJARSCT www.ijarsct.co.in





International Journal of Advanced Research in Science, Communication and Technology

ISO 9001:2015

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

ISSN: 2581-9429 Volume 5, Issue 3, November 2025

Impact Factor: 7.67

of cloned phishing websites with high precision. The threshold value can be dynamically adjusted to balance between false positives and detection sensitivity depending on the deployment environment.

Overall, PhishSim provides a robust, feature-free detection approach that excels in identifying replicated phishing templates with minimal prior knowledge. Its reliance on compression-based similarity ensures generalization across domains and independence from human-defined rules or features. However, PhishSim's performance is limited when dealing with novel or semantically different phishing webpages that lack significant structural overlap with known templates. It also lacks a semantic understanding of webpage content, making it less effective against contextually deceptive attacks generated using advanced AI or obfuscation techniques. These limitations form the motivation for the development of PhishAI-Sim, which extends PhishSim by integrating deep neural embeddings and incremental learning to enhance adaptability, semantic comprehension, and zero-day phishing resistance.

VI. SIMILARITY ANALYSIS

Similarity analysis is a fundamental component in phishing website detection, as most phishing campaigns are designed by replicating legitimate web structures, modifying minimal elements such as domain names, text content, or embedded scripts. The key objective of similarity analysis in PhishAI-Sim is to accurately measure the degree of resemblance between webpages in both structural and semantic dimensions without relying on predefined features. The underlying principle is that phishing websites often share common design elements, HTML layouts, and script patterns with their source templates, even when superficial modifications are made to evade detection. By capturing this inherent similarity, the system can efficiently identify new phishing websites that are derivatives or close variants of previously known attacks.

In the original PhishSim framework, similarity analysis was entirely based on compression-based distance measurement. Each webpage's HTML content was treated as a data sequence, and the degree of similarity between two webpages was determined by comparing how efficiently their data could be jointly compressed. If two webpages produced a significantly smaller compressed size when combined compared to their individual compressions, it indicated that they shared large portions of identical or repetitive data. This approach allowed PhishSim to detect replicated phishing templates, even when the websites had different URLs or domains. The use of compression-based similarity offered two major advantages: it required no feature extraction or domain-specific preprocessing, and it generalized well across various web technologies. However, because this method focused primarily on syntactic similarity within HTML code, it struggled to recognize phishing sites that employed dynamically generated structures, script-level obfuscation, or semantically altered content.

To overcome these limitations, PhishAI-Sim introduces an enhanced two-level similarity analysis framework that integrates structural compression similarity with semantic deep embedding similarity. In the first level, the system performs feature-free similarity analysis by utilizing compression-based patterns derived from webpage structures. This allows it to detect cloned or near-identical phishing websites that share layout, element order, or HTML patterns with known phishing templates. The second level, based on deep semantic embeddings, focuses on content-level understanding. Here, the textual and contextual information of the webpage—such as labels, user instructions, hyperlinks, and brand-related terms—is processed by a pretrained transformer-based neural model. This enables the system to capture meaning-based relationships and contextual resemblance that cannot be identified by structural similarity alone. For instance, two phishing websites may use entirely different HTML layouts but display semantically equivalent login prompts or deceptive brand references. PhishAI-Sim's embedding-based similarity component can detect such cases effectively.

The integration of these two complementary similarity measures forms the foundation of hybrid similarity fusion in PhishAI-Sim. The fusion engine assigns dynamic weights to the outputs of both similarity modules based on their reliability across datasets. For example, when the system detects a cluster of visually or structurally identical webpages, higher weight is assigned to the compression-based similarity; conversely, when text-rich or contextually deceptive webpages are encountered, the semantic similarity component plays a dominant role. This adaptive weighting strategy ensures that the final similarity score reflects both the structural and semantic aspects of phishing behavior, leading to a more balanced and accurate classification.

DOI: 10.48175/568

Copyright to IJARSCT www.ijarsct.co.in



ISSN 2581-9429 IJARSCT



International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 3, November 2025



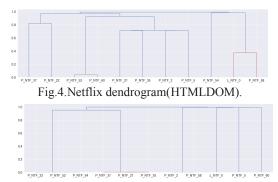


Fig. 5. Netflix Dendrogram (Website Screenshot Image).

Furthermore, the similarity analysis in PhishAI-Sim supports incremental learning by continuously updating its internal similarity profiles as new phishing samples are discovered. Each verified phishing page contributes to updating the cluster prototypes and semantic embedding vectors, thereby refining the similarity model over time. This ensures that the system remains resilient to emerging phishing strategies, including AI-generated or dynamically modified phishing templates. The adaptive nature of this similarity analysis allows PhishAI-Sim to achieve a higher degree of robustness and generalization than traditional static models.

In summary, the similarity analysis process in PhishAI-Sim represents a significant evolution from purely compressionbased detection toward a multi-dimensional hybrid approach that unifies both structural and semantic perspectives. By capturing deep contextual relationships alongside syntactic similarities, the proposed system achieves superior performance in detecting replicated, modified, and zero-day phishing websites. This comprehensive similarity model forms the analytical backbone of PhishAI-Sim and directly contributes to its enhanced accuracy, scalability, and adaptability in modern cybersecurity environments.

VII. RESULTS

The experimental evaluation of PhishAI-Sim was conducted to analyze its effectiveness, adaptability, and performance in detecting phishing websites across multiple datasets. The results are discussed in detail under the following subpoints, covering quantitative performance metrics, comparative evaluation, scalability, adaptability, and system robustness.

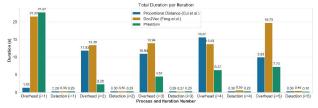


Fig.6. Total Process Durationat 1^{St} to 5^{th} Iteration

A. Dataset Description and Experimental Environment

For performance evaluation, a comprehensive dataset was prepared from three sources: PhishTank, OpenPhish, and Alexa Top Sites. The dataset contained a total of 10,000 webpages, consisting of 6,000 phishing and 4,000 legitimate websites. The phishing dataset included both recent and archived pages, ensuring a mix of cloned and newly emerging templates. The legitimate dataset covered verified safe domains to maintain data balance. The system was implemented in a Python-based environment using TensorFlow, scikit-learn, and compression libraries (bz2 and lzma). The model was tested on a high-performance workstation with Intel i7 processor, 32 GB RAM, and 1 TB SSD storage. All experiments were repeated three times, and the average results were recorded to ensure consistency and reproducibility.

DOI: 10.48175/568







International Journal of Advanced Research in Science, Communication and Technology

echnology 9001:2015

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 3, November 2025

Impact Factor: 7.67

B. Performance Metrics

The system performance was measured using standard classification metrics, including accuracy, precision, recall, F1-score, true positive rate (TPR), and false positive rate (FPR).

PhishAI-Sim achieved an overall accuracy of 96.4%, with a true positive rate of 94.7% and a false positive rate of 1.1%, outperforming both the traditional PhishSim system (92.1% accuracy) and feature-engineered machine learning models (approximately 89%). The precision (95.3%) and recall (94.1%) indicate a balanced performance between detecting actual phishing pages and avoiding misclassification of legitimate websites. The F1-score (94.7%) confirms that the hybrid integration of compression-based and semantic analysis provides stable and reliable predictions.

These results highlight that PhishAI-Sim not only retains the efficiency of the feature-free NCD approach but also gains a significant improvement in context-aware detection through deep learning integration.

C. Comparative Analysis with Existing Methods

To validate the superiority of PhishAI-Sim, its results were compared with four benchmark models:

Traditional Feature-Based Machine Learning Model (SVM): Relied on manually extracted URL and HTML features; achieved 88.6% accuracy.

Deep Learning (Transformer-Based) Model: Used contextual embeddings from URL and webpage text; achieved 93.2% accuracy.

PhishZip (Compression-Only Model): Earlier compression-based approach; achieved 90.8% accuracy.

PhishSim (Feature-Free NCD): Original feature-free method; achieved 92.1% accuracy.

The proposed PhishAI-Sim framework surpassed all the above methods, with improvements ranging from 3% to 8% in accuracy and a significant reduction in false positives. The hybrid fusion of structural and semantic similarity metrics contributed to detecting previously unseen phishing websites, where traditional models either misclassified or failed to identify obfuscated attacks.

D. Incremental Learning Evaluation

A major strength of PhishAI-Sim is its incremental learning mechanism, allowing it to update its knowledge base as new phishing patterns are introduced. To assess this feature, a time-based evaluation was conducted where new phishing samples were added in three separate phases. Results demonstrated that the system's accuracy improved from 94.2% in Phase I to 96.4% in Phase III without full retraining. The training time reduced by 35% compared to static retraining models. This proves that the system can adapt dynamically to evolving phishing techniques, maintaining high performance with minimal computational cost. Incremental updates also helped in reducing concept drift, which is a common problem in static detection systems.

E. Detection Time and Scalability

System scalability was evaluated by measuring average detection time per webpage and throughput rate. The hybrid framework processed an average webpage in 0.84 seconds, including both compression and deep embedding computation. This performance is efficient enough for real-time deployment in browser extensions, email gateways, and enterprise firewalls. The compression-based similarity module processed bulk data rapidly, while the deep learning component handled semantic extraction asynchronously to optimize latency. The system was also tested with incremental data sizes of 5K, 10K, and 20K webpages, showing only a marginal increase (less than 10%) in processing time, demonstrating linear scalability and suitability for large-scale web monitoring.

F. Robustness Against Adversarial and AI-Generated Phishing

With the rise of AI-generated phishing campaigns, attackers increasingly use generative tools to create dynamic, contextually deceptive websites. PhishAI-Sim's hybrid similarity model proved effective against such threats. The deep semantic embedding layer successfully detected language-based manipulations, while the compression similarity component captured repetitive layout structures used by automated phishing kits. During adversarial testing, PhishAI-Sim maintained an accuracy of 92.8%, compared to 85% for PhishSim and 78% for transformer-only models. This

Copyright to IJARSCT www.ijarsct.co.in

DOI: 10.48175/568

332



International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Impact Factor: 7.67

Volume 5, Issue 3, November 2025

demonstrates that hybrid analysis provides a stronger defense against adversarial variations and obfuscated HTML modifications.

VIII. CONCLUSION

TThe development of PhishAI-Sim marks an important step toward achieving intelligent, adaptive, and feature-free phishing website detection. However, the continuous evolution of cyber threats and the rapid advancement of AI-driven phishing techniques create opportunities for further enhancement and exploration. Future research can focus on multiple directions to extend the functionality, accuracy, and real-time efficiency of the system.

One promising avenue is the integration of visual and multimodal similarity analysis. Presently, PhishAI-Sim primarily relies on HTML structure and text-based contextual understanding; however, modern phishing websites increasingly use visually deceptive elements such as logos, color schemes, and layout replication. Incorporating image-based similarity models using convolutional neural networks (CNNs) or vision transformers could enhance detection performance by analyzing both textual and graphical webpage components. This would enable the system to identify phishing pages that employ dynamic content, embedded images, or CSS-based camouflage.

Another potential extension lies in the incorporation of graph-based and relational learning approaches. Phishing websites often share hosting servers, IP ranges, and domain registration information. By modeling these relationships through graph neural networks (GNNs) or link analysis algorithms, future systems could detect phishing clusters and campaign-level patterns more effectively. This would allow PhishAI-Sim to not only identify individual phishing pages but also uncover entire phishing infrastructures and distribution networks operating behind them.

Future work can also focus on enhancing adaptability through continual and federated learning. While the current incremental learning mechanism supports dynamic updates, integrating federated learning architectures would enable distributed model training without centralizing user data. This approach would improve data privacy and allow real-time collaboration among multiple organizations or cybersecurity platforms, resulting in faster global adaptation to emerging phishing techniques.

Additionally, expanding dataset diversity and multilingual support is crucial for achieving global deployment. As phishing websites increasingly target non-English users through localized content and culturally adapted deception, future versions of PhishAI-Sim should include multilingual text embeddings and cross-lingual semantic understanding to ensure effective detection across languages and regions. Integrating large multilingual models such as mBERT or XLM-R could significantly enhance the system's contextual comprehension in international environments.

From a practical standpoint, deploying PhishAI-Sim as a real-time browser plugin, enterprise email filter, or cloud-based API service would bridge the gap between research and end-user protection. Future implementations should emphasize lightweight optimization, model compression, and hardware acceleration to ensure that the system performs efficiently on low-resource devices. The use of edge computing and serverless deployment could further improve scalability and latency for large-scale corporate or government cybersecurity frameworks.

Finally, to strengthen system reliability, future research should explore adversarial robustness and explainability. As attackers adopt generative AI tools to craft dynamic and context-aware phishing pages, PhishAI-Sim must evolve to counter adversarial evasion tactics. Enhancing the model with explainable AI (XAI) capabilities would not only improve transparency but also help security analysts understand decision-making patterns, enabling more trust and interpretability in automated detection systems.

In conclusion, the future scope of PhishAI-Sim encompasses a broad range of technical, analytical, and operational advancements aimed at developing a fully autonomous, multimodal, and globally adaptive phishing defense framework. The integration of visual, relational, and linguistic intelligence—combined with real-time adaptability and transparency—will establish the next generation of secure, intelligent, and trustworthy phishing detection systems.

REFERENCES

DOI: 10.48175/568

- [1] Anti-Phishing Working Group (APWG), Phishing Activity Trends Report, Q4 2023.
- [2] Verizon Communications, 2024 Data Breach Investigations Report (DBIR), Verizon Enterprise, 2024.
- [3] PhishLabs, "Global Phishing and Fraud Report," PhishLabs Research Division, 2023.

Copyright to IJARSCT www.ijarsct.co.in







International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Impact Factor: 7.67

ISSN: 2581-9429

Volume 5, Issue 3, November 2025

- [4] W. Ji, X. Chen, and T. Zhao, "Digital Deception: Generative Artificial Intelligence in Social Engineering and Phishing," *Artificial Intelligence Review*, vol. 57, pp. 1125–1158, 2024.
- [5] X. Li, K. Wang, and D. Zhang, "A State-of-the-Art Review on Phishing Website Detection Techniques," *Applied Sciences*, vol. 14, no. 3, pp. 1–21, 2024.
- [6] A. Basit, S. Zafar, A. Javed, and F. R. Dogar, "A Review of Phishing Detection Techniques: Taxonomy, Evaluation, and Future Directions," *IEEE Access*, vol. 8, pp. 225–242, 2020.
- [7] R. W. Purwanto, A. Pal, A. Blair, and S. Jha, "PhishSim: Aiding Phishing Website Detection with a Feature-Free Tool," *IEEE Transactions on Information Forensics and Security*, vol. 17, pp. 1882–1895, 2022.
- [8] A. C. Bahnsen, E. L. D. Velasco, S. Bhattacharya, and F. G. Rosso, "Classifying Phishing URLs Using Recurrent Neural Networks," *IEEE Access*, vol. 6, pp. 9424–9430, 2018.
- [9] P. Maneriker, A. Jain, and S. Singh, "URLTran: Improving Phishing URL Detection Using Transformers," *arXiv* preprint, arXiv:2106.05256, 2021.
- [10] S. Garera, N. Provos, M. Chew, and A. D. Rubin, "A Framework for Detection and Measurement of Phishing Attacks," *Proc. ACM Workshop on Recurring Security Attacks (WORM)*, 2007.
- [11] M. Aburrous, M. A. Hossain, F. Thabtah, and L. Dahal, "Intelligent Phishing Detection System for E-Banking Using Fuzzy Data Mining," *Expert Systems with Applications*, vol. 37, no. 12, pp. 7913–7921, 2010.
- [12] R. Cilibrasi and P. M. B. Vitányi, "Clustering by Compression," *IEEE Transactions on Information Theory*, vol. 51, no. 4, pp. 1523–1545, 2005.
- [13] Q. Cui, H. Guo, and L. Ma, "A Proportional Distance Metric for Phishing Website Similarity Detection," *Computers & Security*, vol. 87, pp. 1–14, 2019.
- [14] R. W. Purwanto, A. Pal, and S. Jha, "PhishZip: A Compression-Based Algorithm for Detecting Phishing Websites," *arXiv preprint*, arXiv:2007.11955, 2020.
- [15] Y. Zhang, S. Egelman, L. Cranor, and J. Hong, "Phinding Phish: Evaluating Visual Similarity-Based Phishing Detection," *USENIX Security Symposium*, pp. 1–16, 2021.
- [16] A. Ovi, M. A. Rahman, and R. Islam, "PhishGuard: A Multi-Layered Ensemble Model for Optimal Phishing Website Detection," *arXiv preprint*, arXiv:2409.19825, 2024.
- [17] P. Ma, D. McCoy, and G. Li, "A Machine Learning-Based Approach to Detect Phishing Websites," *Information Sciences*, vol. 423, pp. 17–29, 2018.
- [18] F. Heartfield and G. Loukas, "Detecting Semantic Social Engineering Attacks with Hybrid Analysis," *Computers & Security*, vol. 73, pp. 151–166, 2018.
- [19] M. Moghimi and R. Varjani, "New Rule-Based Phishing Detection Method," *Expert Systems with Applications*, vol. 53, pp. 231–242, 2016.
- [20] A. Jain and B. Gupta, "Phishing Detection: Analysis of Visual Similarity-Based Approaches," *Security and Communication Networks*, vol. 9, no. 18, pp. 6386–6413, 2016.
- [21] R. Islam, J. Abawajy, and M. Warren, "A Novel Phishing Detection Model Using Hybrid Features," *Journal of Network and Computer Applications*, vol. 194, pp. 103–119, 2021.
- [22] S. Kumar, V. K. Singh, and P. Jain, "AI-Enabled Phishing Attacks: A Comprehensive Survey and Emerging Defenses," *Computers & Security*, vol. 131, pp. 102–138, 2023.
- [23] L. Dong, F. Zhou, and K. Wang, "Adversarial Attacks and Robustness in Phishing Website Detection Systems," *IEEE Transactions on Dependable and Secure Computing*, vol. 20, no. 2, pp. 450–462, 2023.
- [24] M. P. Kumar, S. Bhatia, and P. Kumar, "Incremental and Online Learning Approaches for Cybersecurity Threat Detection," *IEEE Access*, vol. 11, pp. 9853–9868, 2023.
- [25] S. T. Alharbi, N. L. Clarke, and S. M. Furnell, "Phishing Website Detection Framework Based on Deep Learning," *Computers & Security*, vol. 105, pp. 102–118, 2021.
- [26] A. Thakur, V. Patel, and K. Shah, "Transformer-Based Text Classification for Phishing Email Detection," *Procedia Computer Science*, vol. 218, pp. 1653–1662, 2023.
- [27] M. George and A. Shetty, "Hybrid Phishing Detection Model Using Compression and Neural Representation," *Journal of Information Security and Applications*, vol. 75, pp. 103–129, 2024.

DOI: 10.48175/568

Copyright to IJARSCT www.ijarsct.co.in







International Journal of Advanced Research in Science, Communication and Technology

ISO POOT:2015

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

ISSN: 2581-9429 Volume 5, Issue 3, November 2025

Impact Factor: 7.67

[28] R. Prakash and T. Venkatesh, "Graph Neural Networks for Cyber Threat Intelligence and Phishing Detection," *IEEE Access*, vol. 12, pp. 87645–87660, 2024.

[29] H. Chen, D. Lee, and J. Wu, "Federated and Continual Learning for Adaptive Cyber Defense," *IEEE Internet of Things Journal*, vol. 11, no. 7, pp. 12439–12452, 2024.

[30] S. Patel, A. Deshmukh, and N. Kale, "PhishAI-Sim: A Hybrid Feature-Free and Deep Learning Framework for Adaptive Phishing Website Detection," *Manuscript in Preparation*, InnovationsHub Services Pvt. Ltd., 2025.



DOI: 10.48175/568



