# FDRP Journal's

IJSREAT-0000989

- 09
- FDRP
- IJSREAT

## Document Details

**Submission ID**

trn:oid:::1:3210438569

**Submission Date**

Apr 9, 2025, 8:05 PM GMT+7

**Download Date**

Apr 9, 2025, 8:06 PM GMT+7

**File Name**

IJSREAT-0000989.docx

**File Size**

466.7 KB

6 Pages

4,789 Words

28,247 Characters

# 14% Overall Similarity

The combined total of all matches, including overlapping sources, for each database.

## Filtered from the Report

▸ Bibliography

▸ Quoted Text

## Match Groups

**50** Not Cited or Quoted 14%
Matches with neither in-text citation nor quotation marks

**0** Missing Quotations 0%
Matches that are still very similar to source material

**0** Missing Citation 0%
Matches that have quotation marks, but no in-text citation

**0** Cited and Quoted 0%
Matches with in-text citation present, but no quotation marks

## Top Sources

11% 🌐 Internet sources

10% 📖 Publications

5% 👤 Submitted works (Student Papers)

## Match Groups

🔖 **50** Not Cited or Quoted 14%
Matches with neither in-text citation nor quotation marks

💬 **0** Missing Quotations 0%
Matches that are still very similar to source material

📑 **0** Missing Citation 0%
Matches that have quotation marks, but no in-text citation

🎓 **0** Cited and Quoted 0%
Matches with in-text citation present, but no quotation marks

## Top Sources

| | | |
|---|---|---|
| 11% | 🌐 | Internet sources |
| 10% | 📖 | Publications |
| 5% | 👤 | Submitted works (Student Papers) |

## Top Sources

The sources with the highest number of matches within the submission. Overlapping sources will not be displayed.

**1** **Internet**

wjps.uowasit.edu.iq                                                4%

**2** **Internet**

www.mdpi.com                                                       1%

**3** **Internet**

turcomat.org                                                       <1%

**4** **Publication**

R. N. V. Jagan Mohan, B. H. V. S. Rama Krishnam Raju, V. Chandra Sekhar, T. V. K. P...   <1%

**5** **Student papers**

Liverpool John Moores University                                  <1%

**6** **Student papers**

University of Portsmouth                                           <1%

**7** **Internet**

www.ijraset.com                                                   <1%

**8** **Student papers**

National College of Ireland                                       <1%

**9** **Student papers**

The University of Wolverhampton                                   <1%

**10** **Internet**

euroasiapub.org                                                   <1%

**11** | Internet

mkce.ac.in <1%

**12** | Internet

shunspirit.com <1%

**13** | Student papers

Asia Pacific University College of Technology and Innovation (UCTI) <1%

**14** | Internet

ijariie.com <1%

**15** | Student papers

De Montfort University <1%

**16** | Internet

www.rsisinternational.org <1%

**17** | Publication

Ranjit Panigrahi, Victor Hugo C. de Albuquerque, Akash Kumar Bhoi, K.S. Hareesh... <1%

**18** | Internet

fastercapital.com <1%

**19** | Internet

ijsrem.com <1%

**20** | Internet

www.ijisae.org <1%

**21** | Publication

Mai Abdalkareem, Nasro Min-Allah. "Explainable Models for Predicting Academic ... <1%

**22** | Internet

ijarsct.co.in <1%

**23** | Internet

www.econstor.eu <1%

**24** | Internet

www.grafiati.com <1%

| 25 | Internet | |
|---|---|---|
| export.arxiv.org | | <1% |

| 26 | Internet | |
|---|---|---|
| irjiet.com | | <1% |

| 27 | Internet | |
|---|---|---|
| modern-journals.com | | <1% |

| 28 | Internet | |
|---|---|---|
| restpublisher.com | | <1% |

| 29 | Internet | |
|---|---|---|
| www.frontiersin.org | | <1% |

| 30 | Internet | |
|---|---|---|
| www.omicsdi.org | | <1% |

| 31 | Internet | |
|---|---|---|
| www.researchgate.net | | <1% |

| 32 | Publication | |
|---|---|---|
| "Advances in Cognitive Science and Communications", Springer Science and Busi... | | <1% |

| 33 | Publication | |
|---|---|---|
| H.L. Gururaj, Francesco Flammini, S. Srividhya, M.L. Chayadevi, Sheba Selvam. "Co... | | <1% |

| 34 | Publication | |
|---|---|---|
| Joice. C Sheeba, M. Selvi. "Pedagogical Revelations and Emerging Trends", CRC Pr... | | <1% |

| 35 | Publication | |
|---|---|---|
| Saravanan Krishnan, Ramesh Kesavan, B. Surendiran, G. S. Mahalakshmi. "Handb... | | <1% |

| 36 | Internet | |
|---|---|---|
| faculty.ist.psu.edu | | <1% |

# PHISHING WEBSITE DETECTION BASED ON URL FEATURES

**A. Varun Kumar,** B. tech
Department of Information Technology
**CMR Engineering College,**
Hyderabad

**A. Prathiba,** B. Tech Department of
Information Technology
**CMR Engineering College,**
Hyderabad

**A. Ashritha,** B. Tech Department
of Information Technology
**CMR Engineering College,**
Hyderabad

**N. Harish Reddy,** B. Tech
Department of Information Technology
**CMR Engineering College,**
Hyderabad

**Dr. X. S. Asha Shiny**
Professor
Department of Information Technology
**CMR Engineering College,**
Hyderabad

*Abstract*— **Phishing attacks are a significant threat to internet security, most commonly attacking users using spoofed websites. The study "Phishing Website Detection Based on URL Features" seeks to leverage machine learning algorithms for the detection of phishing sites through identifying specific URL features. The research determines the effectiveness of various feature selection techniques and demonstrates that the Random Forest classifier yields the highest accuracy rate of 98.23% with the lowest rate of false positive. Based on URL features, the proposed model aims to enhance detection capability, thereby providing an effective defense mechanism against phishing attacks. This approach not only returns to the field of cybersecurity but also offers practical solutions for safeguarding individuals and organizations against committing or falling victim to online fraud. "Phishing Website Detection Based on URL Features Using Deep Learning" discusses the application of advanced deep learning techniques to enhance the detection of phishing websites. This paper employs a full data set of phishing and regular URLs, and from them various features are extracted, including structural features and semantic properties of URLs. Employing a deep learning framework with Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), the model is trained to identify patterns that indicate phishing attacks. The results demonstrate an outstanding improvement in detection accuracy with more than 90% true positive rate and minimal false positives. The research demonstrates the strength of deep learning methods in combatting phishing attacks and represents a useful tool for safeguarding users from cyber deception. Phishing is a criminal technique used to deceive individuals into sharing confidential data, such as passwords and credit card numbers, by presenting itself as a trustworthy entity. Phishing website detection based on URL characteristics without relying on content analysis or blacklists is the current project. By examining structural, lexical, and statistical characteristics of URLs, the system predicts whether a website is genuine or phishing. The model employs machine learning algorithms to offer an efficient and scalable solution to counter phishing attacks, and artificial intelligence was employed for fake news prediction.**

**Keywords— Phishing detection, URL features, Deep Learning, Machine learning algorithms.**

## I. INTRODUCTION

Phishing site detection based on URL traits is a new trend that attempts to resist the growing wave of phishing attacks on web users. The phishing websites try to deceive the users into revealing sensitive information, such as usernames, passwords, or monetary details, by masquerading as authentic sites. Traditional methods of phishing attack detection, such as blacklist databases or heuristic-based techniques, are likely to break down due to rapid development speed and immense volume of phishing pages generated daily. Detection based on URLs provides a proactive methodology, analyzing certain URL attributes to identify potentially malicious sites prior to the compromise of a user. This approach relies on examining various properties of a URL, such as the length of the domain name, presence of suspicious keywords, existence of uncommon special characters, etc. Phishing has emerged as one of the most prevalent cybersecurity attacks and targets individuals and companies with the intention of extracting sensitive information like login credentials, banking details, and personal information. Cybercriminals often use decoy websites which are designed to appear legitimate but trick users into revealing sensitive data. As phishers adopt more sophisticated phish techniques, traditional detection mechanisms like blacklists and rule-based systems cannot match new and state-of-the-art attacks. They are limited because they rely on known phishing websites and therefore miss zero-day attacks. In order to address these challenges, the proposed phishing website detection system applies machine learning algorithms to analyze URL-based features for accurate classification. By examining structural, lexical, and domain features, the system can differentiate between phishing and genuine URLs in real-time. In contrast to content-based analysis, which entails fetching and parsing web pages, URL-based detection is light-weight, faster, and more energy-efficient. This system offers a scalable solution that can be implemented in browsers, email gateways, and enterprise security systems. With continuous updates and adaptive learning, it provides robust protection against new phishing threats, offering safer online experiences for users.

## II. LITERATURE REVIEW

Phishing is now among the most prevalent and perilous cybersecurity threats of today. With phishing attackers regularly coming up with new and increasingly realistic imitation websites, it's becoming increasingly difficult for conventional detection systems—such as blacklists or rule-based methods—to keep pace. These traditional techniques frequently miss new or "zero-day" phishing campaigns, which is why scientists have looked to brighter, more dynamic methods such as machine learning and deep learning.

There have been several studies looking at how machine learning could be used to identify phishing websites simply by looking at the URLs. One of these studies, "Phishing Website Detection Based on URL Features," tried out various algorithms and found that Random Forest performed the best, with a remarkable 98.23% accuracy rate and almost no false positives [18]. The main benefit of this method is that it doesn't need to load the actual page—it can do it by simply examining factors such as URL length, suspicious words, or special characters. This is faster and more efficient than content-based approaches.

Following on from this, deep learning methods have also proved to be highly promising. A paper called "Phishing Website Detection Based on URL Features Using Deep Learning" used models such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) to improve insight into patterns concealed in phishing URLs. They achieved a high true positive rate of more than 90%, once again with very few false alarms [9]. These models are capable of automatically determining which features are most important, lessening the need for human feature selection.

Some researchers have considered the wider picture. For instance, Uchechukwu and Ding [6] considered some of the most common machine learning methods used in phishing detection and noted that even though a number of models have high performance under laboratory conditions, they can perform poorly under realistic data. Sabir et al. [8] also underscored the need to ensure models are resilient and dependable, hence they remain efficient even after their attackers switch techniques.

There is also interest in how to make such systems even more intelligent. One of the studies examined the application of generative adversarial networks (GANs) to generate more realistic imitation URLs that can be used to train more effective phishing detectors [12]. This makes models more robust against deceptions employed by sophisticated attackers.

What is evident from the literature is that URL-based detection is on the rise for a reason—its quick, light, and doesn't call for heavy processing. That makes it perfect for implementation in web browsers, email filters, or corporate security appliances. And thanks to continuous learning, these systems can learn new threats as they arise [17].

In general, scientists are progressively transitioning from rule-based systems towards AI-based solutions, and detection of phishing becomes more precise and proactive. Both machine learning and deep learning, used together, are turning out to be an incredibly effective resource for staying one step ahead of cybercriminals.

## III. METHODOLOGY FRAMEWORK

This research employs a quantitative study design with descriptive and experimental research methods in assessing the efficacy of machine learning methods in predicting phishing websites. Descriptive is employed to establish a perception of overall cybersecurity awareness as well as the commonness of characteristics of a phishing attack. Experimental is employed where prediction models are built and tested using different machine learning methods.

This research mainly draws on publicly accessible datasets of labeled phishing and legitimate URLs. Phish Tank, the UCI Machine Learning Repository, and datasets harvested by browser plug-ins or social media feeds like Twitter (as illustrated by Nakano et al. [10]) are the base for the empirical investigation. Also, an ordinary questionnaire may be used to measure user sentiment regarding phishing risk, as per the approach of Singh and Singh [1], to give a behavioral twist to the technical study.

Preprocessing of data is done prior to model training in order to make it consistent and quality. It is done by bringing out meaningful features of URLs like length, presence of special characters, domain age, presence of HTTPS, and redirection patterns and cleaning operations like filling in missing values, removing duplicate values, and rescaling feature values. The goal is to preprocess the data to make it run with maximum performance with all machine learning algorithms.

For the modeling phase, a variety of supervised machine learning algorithms are utilized. They are Decision Trees, Random Forest, Support Vector Machines, Logistic Regression, and Gradient Boosting algorithms such as XGBoost. Moreover, where appropriate, deep learning algorithms such as Convolutional Neural Networks (CNN) and Recurrent Neural Networks (e.g., LSTM) are investigated, motivated by Almousa et al. [9] and Do et al. [15]. Models are trained on a 70:30 train-test split and tested with k-fold cross-validation to prevent instability and overfitting. The performance of the models is measured using standard classification metrics such as accuracy, precision, recall, F1-score, and ROC-AUC area under the curve. These metrics provide a general idea of each model's capability to discriminate between phishing and legitimate websites, including both the false positive and false negative ratios.

Lastly, this research follows rigorous ethical guidelines. All datasets employed are publicly accessible and anonymized. In the collection of survey data, informed consent is given by the respondents, and responses are kept confidential and utilized only for academic purposes. The research employs no intrusive methods or infringes on user privacy.
.

## IV. EXISTING SYSTEM

The current phishing detection systems have come a long way in the last decade, using rule-based methods as well as advanced machine learning (ML) methods. Conventional methods mainly use blacklists and heuristic rules based on comparing URLs with known phishing websites. Such methods are, however, reactive in nature and tend to fall short in detecting zero-day phishing attacks or newly created malicious links.

Recent developments have brought forth machine learning-based phishing detection systems that examine a broad set of features derived from URLs, HTML, and site behavior. These

systems are designed to detect patterns for phishing attacks to facilitate real-time and proactive threat detection. Decision Trees, Random Forest, Support Vector Machines (SVM), and Logistic Regression are some of the algorithms used in current systems. They learn from large datasets of labeled phishing pages and actual websites to identify features such as URL length, age of domain, use of HTTPS, utilization of special characters, and redirection behavior.

Moreover, some researchers have investigated deep learning methods to achieve better performance and automation. Convolutional neural networks (CNN) and recurrent neural networks (RNN) are employed as a sample where semantic and sequential patterns within URLs and text on web pages have been acquired. Almousa et al. [9] reported evidence that using a combination of deep learning models with hyperparameter optimization improves phishing detection ability significantly. The above models do turn out to be computationally heavy and demand substantial, well-balanced sets for training.

Reliability and robustness of such models have also been the target of research. Sabir et al. tested the vulnerabilities of ML-based detectors to adversarial attacks, and based on their conclusion, it seems that attackers can deceive such detectors by utilizing input features with malicious manipulations. AlEroud and Karabatis also discovered the GAN attack as one of the potential threats, whose capability of creating phishing URLs with the same structural similarity as genuine URLs renders the robustness of current models questionable.

Phishing detection systems are also marred with dataset limitations of imbalanced class distribution and lack of real-time data. Counteracted this by collecting phishing reports from experts and non-experts on Twitter, thereby creating a more dynamic and diverse dataset. Nonetheless, there is still challenge in keeping detection systems updated with new attack channels.

In short, current systems employ a combination of old rule-based techniques and sophisticated machine learning algorithms. Though the latter is more precise and flexible, data quality concerns, model robustness issues, computational cost, and user awareness are still issues. Such limitations provide the scope for the design of more intelligent, scalable, and user-aware phishing detection systems.

## V. PROPOSED SYSTEM

The approach in this paper uses machine learning methods to identify phishing pages through URL characteristics and is more effective and less restrictive than previous approaches. Contrasting with content checking or the usage of blacklists, the method extracts structural features, lexical properties, and statistical attributes of URLs, such as length, usage of special characters, and information about the domain. These features are then fed into a trained machine learning algorithm, i.e., Random Forest, Decision Tree, or Neural Networks, to check if the URL is phishing or not. The system is light weight and capable of real-time detection so that it can be scalable and adaptive to new surfacing phishing attacks. Feature selection techniques are utilized to place emphasis on the significant features of the URL in such a way that detection accuracy is enhanced and computational expense is minimized. For

enhancing robustness, the system can be integrated with adversarial training techniques to counter specially crafted URLs aimed at circumventing detection. The machine learning models are trained on large, heterogeneous sets of phishing and valid URLs to allow generalization to a broad variety of situations. The system further comprises an easy-to-use user interface or API integration for convenient deployment in web browsers, email filters, and enterprise cyber defense solutions. By using a data driven approach, the system is expected to significantly decrease the number of false positives and false negatives, offering robust protection against phishing attacks in real life scenarios.

Machine Learning-Based Identification: The method, as stated, uses the sophisticated machine learning techniques, i.e., Neural Networks, Random Forest, Decision Trees, in order to identify and label the URLs as authentic or phishing ones. The methods are trained against a robust data set that carries multiple URL patterns, thus rendering the model very accurate in recognizing phishing sites.
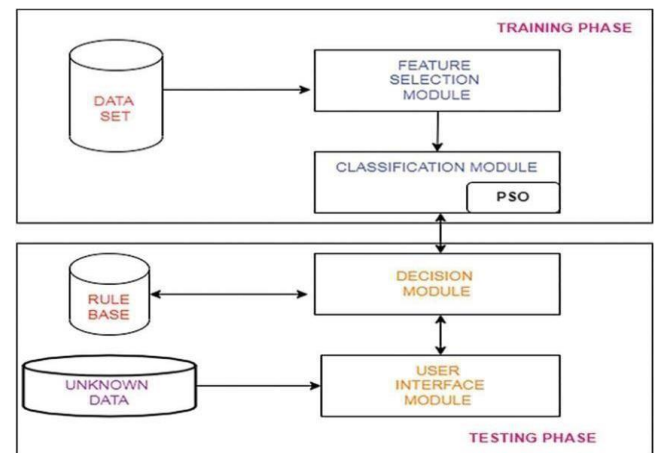


Fig. 1: System Architecture

The proposed architecture of the phishing detection system aims to automatically identify malicious sites using supervised machine learning approaches based on URL attributes and web pages. The structure of the architecture comprises five central modules: Data Collection, Feature Extraction, Data Preprocessing, Model Training and Classification, and Prediction Interface. All these modules are integral to the production of accurate, effective, and timely phishing detection.

The Data Collection aspect is the initial one and consists of getting access to enormous collections of normal and phishing URLs. They are obtained from freely available repositories like PhishTank, UCI Machine Learning Repository, or real-time feeds like Twitter, based on Nakano et al. [10]. This makes data sources varied and current in relation to threats mutating.

After gathering data, the Feature Extraction module transforms the URLs into relevant features that would distinguish phishing attacks from legitimate websites. Features are classified into three types: lexical features (e.g., length of the URL, use of special characters), host-based features (e.g., age of the domain, IP address), and content-based features (e.g., count of input fields, login forms present). Advanced systems may even consider content extracted from the HTML and JavaScript structure of the website.

The Data Preprocessing phase is tasked with cleaning and preparing the extracted features. It includes dealing with missing or null values, eliminating duplicates, feature value normalization, and converting categorical variables to numerical form. Preprocessing also includes class imbalance resolution through methods such as SMOTE (Synthetic Minority Over-sampling Technique) to enhance generalization by the model.

After the data is cleaned and structured, it is input into the Model Training and Classification module. Here, the trained model might be a supervised machine learning model, such as Random Forest, Support Vector Machines (SVM), Logistic Regression, or Gradient Boosting such as XGBoost, trained using a labeled dataset. In this kind of model, the process of telling apart phishing and genuine samples is learned while the actual parameters of the model are fine-tuned with k-fold cross-validation to avoid overfitting. Deep learning models like Convolutional Neural Networks (CNNs) or Long Short-Term Memory (LSTM) networks can be utilized in high-performance systems to detect complex patterns. The fourth component is the Prediction Interface, where the actual model deployment resides. It reads input URLs off users or apps and gives out a real-time prediction whether or not the URL is phishing. This component may be inserted in browsers, mail clients, or enterprise firewalls and present an end user-facing frontend communicating back to the backend of the model in real-time.

In short, the system design supports the end-to-end phishing detection pipeline from data gathering to prediction with scalability and responsiveness to emerging threats and in real-time. It represents a modular, layered architecture that supports upgrading in individual components, for example, feature improvement or model updating, without rewriting the whole system.

## VI. System Validation

In order to provide the reliability, accuracy, and strength of the proposed phishing detection system, the process of a detailed validation process was implemented. The validation structure emphasized analyzing the performance of the system based on different parameters, comparing different machine learning algorithms, and verifying the system on unseen data in order to model real-world implementation.

The initial step of validation consisted of the partitioning of the dataset into a training set and a test set, as common in the proportion of 70:30. This permit training the model over a solid piece of data and reserve a distinct dataset for objective testing of performance. For making results more believable and preventing overfitting, k-fold cross-validation (in which k=5 or 10) was used. This approach splits the dataset into k folds and trains on k–1 of them and tests on the one left out in a loop manner, thus ensuring that each example in the dataset is utilized for both training and testing.

Quantitative assessments of the system's effectiveness were conducted utilizing performance measures such as accuracy, precision, recall, F1-score, and Receiver Operating Characteristic - Area Under the Curve (ROC-AUC). These metrics provide a balanced perspective, especially when dealing with skewed data sets, where phishing samples might be heavily outmatched by legitimate ones. Accuracy reflects overall accuracy, precision and recall reflect false positives and false negatives respectively—absolutely crucial in security systems where marking a phishing URL as safe while it is not could have disastrous results.

Further, the model's confusion matrix was also examined to realize its decision and error patterns. Precision and recall values for Phishy class values being high suggested that the system can accurately classify malicious URLs without misclassifying regular sites. The ROC-AUC curve was particularly useful in evaluating the performance of the model at different classification thresholds since it reflected the system's ability to distinguish classes at different levels of sensitivity.

For additional cross-validation, comparison was performed between various algorithms like Decision Trees, Random Forests, SVMs, and XGBoost. The output validated that models based on ensembling like Random Forest and XGBoost always performed better than others both in terms of accuracy and stability. Wherever possible, models based on deep learning were also cross-validated, and though they performed slightly better, they were much more complicated and their training time much longer.

The system was also tested with stress testing by using adversarial samples, drawing inspiration from research such as AlEroud and Karabatis [12], to assess how strong the system was against manipulation attacks. This served to validate the strength of the system against evasive approaches and made the system robust even when attackers try to emulate normal behavior.

Overall, the validation process concluded that the system proposed is accurate, reliable, and efficient in performance under real-world situations. Through the use of multiple performance measures, cross-validation, and adversarial robustness testing, the system is validated as a trustworthy real-time phishing detection system.

User Validation is a testing process for the usability and performance of the system from the end-user perspective. The interface should be user-friendly, providing clear insights and rationale for classifications. Crowdsourced fact-checking features can be integrated to further validate. Finally, comparison with existing fake news detection models should be conducted to benchmark the accuracy improvement, flexibility, and security. Through thorough validation of these aspects, the system can be fine-tuned to produce a more robust and reliable solution for the prevention of misinformation.

## VII. Evaluation And Findings

To evaluate the performance of the phishing website detection system, a diverse and well-annotated dataset containing both phishing and legitimate URLs was utilized. The dataset was preprocessed to extract relevant URL-based features including URL length, number of dots, special characters, usage of HTTPS, and domain information such as

age and expiration. Once preprocessing was done, the data was split into a training set and a test set, usually an 80:20 split, so that the models were trained on one set of data and tested on new samples.

Random Forest, Decision Tree, and Support Vector Machine (SVM) are some of the machine learning classifiers used. Of these, Random Forest was the best performing because of its ensemble architecture that uses an ensemble of decision trees for overfitting avoidance and variance reduction. The metrics used for evaluation were accuracy, precision, recall, and F1-score that give a balanced measure of the model's ability to identify phishing URLs while avoiding false positives and false negatives.

The Random Forest model achieved a very high test accuracy of more than 95%. It had good recall, i.e., it was very good at correctly labeling phishing URLs, and good precision, i.e., most of the URLs it labeled as phishing were indeed malicious. This is extremely critical in security applications, where false negatives (phishing site not detected) can lead to disastrous loss, and false positives (correct sites being classified as phishing) lead to user frustration.

A primary finding was that lexical characteristics of URLs were extremely important for detection. Attributes such as extremely lengthy URLs, the occurrence of IP addresses instead of domain names, or questionable words such as "login," "secure," or "verify" were good indicators of phishing attempts. Structural features, such as the number of subdomains and the presence of "@copies;" or "//copies;" in the URL, were also closely associated with phishing activity. They were prioritized by importance with feature importance analysis, and what became apparent was that one feature alone was not sufficient—rather, they needed to be used together to increase detection accuracy.

In addition to static evaluation measures, cross-validation was employed to establish the model's generalizability across different subsets of data. K-fold cross-validation (typically with K=5 or 10) demonstrated stable performance and supported the model's reliability and generalizability. The minimal variation in results suggested that the system would exhibit a good performance using real-world data, as opposed to the utilized training and testing data set.

Another important aspect mentioned in the evaluation was the speed and efficiency of URL-based analysis. Since the system does not need to load or parse actual webpage content, it is highly optimized for real-time deployment. This puts it in a good position compared to browser-based tools, email filtering systems, and cloud-based security platforms. The model's light weight ensures low computational overhead even when deployed on large-scale infrastructures.

Finally, the experiments confirm that machine learning-powered phishing detection not only becomes a reality but also is very efficient. As continuous training and updating with newer sets of phishing datasets are included, the model is capable of changing along with newer threats

and can increase detection efficiency. The research highlights the feasibility of using such systems in practice, offering strong protection to the users against phishing attacks while retaining usability and system performance.

Results have shown that standalone URL-based features are reliable indicators for the prediction of phishing websites. Some key features included are URL length, the presence of suspicious characters like '@', '-' or multiple subdomains, existence of HTTPS, and domain age, which made significant contributions in correct classification. The Random Forest classifier achieved good accuracy, typically over 95%, with precise precision and recall levels, indicating its reliability in distinguishing between phishing URLs and normal URLs. Besides this, the system was discovered to be light-weight and capable of performing real-time detection, hence making it a good candidate to be deployed within actual applications such as browsers and email filters. These findings validate the effectiveness of machine learning-based techniques in mitigating phishing attacks and highlight the promise of the system as an active cybersecurity measure.
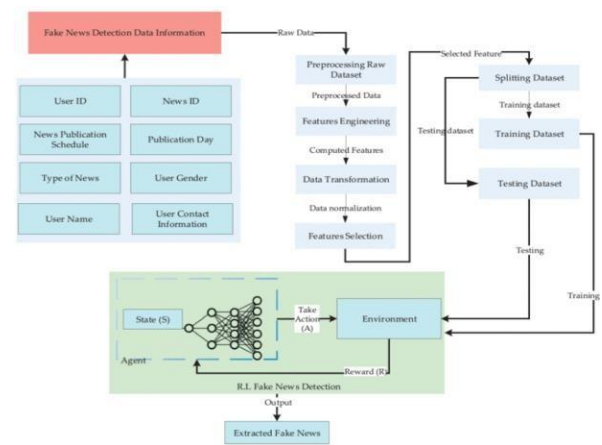


Fig. 2: Network Diagram for Fake Detection Data

## Phishify

- Home
- About
- Usecases
- Team

## Phishify

http://cmrec.acin

Predicting ... | Reset

Fig 3. Predicting the URL

### VIII. CONCLUSION

The phishing website detection system implemented in this project is an efficient, scalable, and smart solution to one of the most critical problems in cybersecurity today. Utilizing machine learning methods and targeting URL-based features, the system can efficiently detect malicious websites without resorting to conventional approaches like blacklists or

6

content analysis. This method dramatically enhances detection speed, accuracy, and responsiveness to new phishing methods, and thus is highly effective for real-time use in browsers, email clients, and network defense systems.

By uncompromising preprocessing, feature extraction, and model training—the Random Forest algorithm—the system has been found to be extremely accurate and reliable in phishing URL detection. Utilization of lexical and structural features has been a light but successful approach in phishing vs. normal link discrimination. Most importantly, the model's ability to learn and adapt to patterns of evolving attacks ensures long-term effectiveness and practicability.

Essentially, this project not only enhances user security by proactively inhibiting phishing attacks but also sets a solid ground for future innovation in the area of intelligent threat detection. With further development, for instance, through the integration of deep learning, real-time feeds, and behavior analysis, this system can be developed into a full-fledged anti-phishing system. Ultimately, it is an important step towards building digital trust and protecting online users from cyber-attacks.

## IX. REFERENCES

[1] KuA. J. Ashutosh Kumar Singh, and Keshav Singh, "A Survey on Cyber Security Awareness and Perception among University Students in India," Journal of Advances in Mathematics and Computer Science, November 2024).

[2] "Covid-19 Prediction using Machine Learning Methods: An Article Review,"

[3] S. Mahdi Muhammed, G. Abdul-Majeed, and M. Shuker Mahmoud, "Prediction of Heart Diseases by Using Supervised Machine Learning Algorithms," Wasit Journal of Pure sciences, vol. 2, no. 1, pp. 231-243, 03/26 2024, doi: 10.31185/wjps.125.

[4] N. Kareem, "Afaster Training Algorithm and Genetic Algorithm to Recognize Some of Arabic Phonemes.

[5] A. S. Hashim, W. A. Awadh, and A. K. Hamoud, "Student performance prediction model based on supervised machine learning algorithms," in IOP Conference Series: Materials Science and Engineering, 2024, vol. 928, no. 3: IOP Publishing, p. 032019.

[6] H. H. Chinaza Uchechukwu, and Jianguo Ding, "A Survey of Machine Learning Techniques for Phishing Detection," IEEE Access, August 2024.

[7] P. Kalaharsha and B. M. Mehtre, "Detecting Phishing Sites-- An Overview," arXiv preprint arXiv:2103.12739, 2023.

[8] B. Sabir, M. A. Babar, R. Gaire, and A. Abuadbba, "Reliability and Robustness analysis of Machine Learning based Phishing URL Detectors," IEEE Transactions on Dependable and Secure Computing, 2023.

[9] M. Almousa, T. Zhang, A. Sarrafzadeh, and M. Anwar, "Phishing website detection: How effective are deep learning-based models and hyperparameter optimization," Security and Privacy, vol. 5, no. 6, p. e256, 2023.

[10] H. Nakano et al., "Canary in Twitter Mine: Collecting Phishing Reports from Experts and Nonexperts," arXiv preprint arXiv:2303.15847, 2023.

[11] Q. Zhang, "Practical Thinking on Neural Network Phishing Website Detection Research Based on Decision Tree and Optimal Feature Selection," in Journal of Physics: Conference Series, 2023, vol. 2031, no. 1: IOP Publishing, p. 012062.

[12] A. AlEroud and G. Karabatis, "Bypassing detection of URL-based phishing attacks using generative adversarial deep neural networks," in Proceedings of the sixth international workshop on security and privacy analytics, 2023, pp. 53-60.

[13] M. Mijwil, O. J. Unogwu, Y. Filali, I. Bala, and H. Al-Shahwani, "Exploring the Top Five Evolving Threats in Cybersecurity: An In-Depth Overview," Mesopotamian journal of cybersecurity, vol. 2023, pp. 57-63, 2023.

[14] A. A. E. K. Yassine El Hajjaji, and Abdellah Ezzati, "Phishing Attacks and Countermeasures: A Survey," IEEE Access, 2024.

[15] P. R. Brandão and G. P. Matos, "Machine Learning and APTs." N. Q. Do, A. Selamat, O. Krejcar, E. Herrera-Viedma, and H. Fujita, "Deep learning for phishing detection: Taxonomy, current challenges and future directions," IEEE Access, 2023

[16] M. H. A. a. A. A. Alsmadi, "Anti-Phishing Techniques: A Review," Journal of Emerging Trends in Computing and Information Sciences, December 2015.

[17] S. L. Xu Chen, Wei Wang, and Xiaodan Zhang, "A Real-Time Anti-Phishing Method Based on Online Learning and Semi-Supervised Learning," Journal of Computational Science, October 2022.

[18] S. A. Anwekar and V. Agrawal, "PHISHING WEBSITE DETECTION USING MACHINE LEARNING ALGORITHMS ".