

Suspicious URL Checker with Cloud ML

Vinusha S, Dr .E.Manohar

Computer Science, Francis Xavier Engineering College, Tirunelveli – TamilNadu - India
vinushas.ug22.cs@francisxavier.ac.in

Computer Science, Francis Xavier Engineering College, Tirunelveli – TamilNadu-India
manohar@francisxavier.ac.in

Abstract:

People are online all the time now, and that just makes things easier for cybercriminals. Phishing emails, dodgy websites, or a homepage that suddenly looks like it was hacked—they're everywhere. What really gets under most people's skin is how these attacks hide behind everyday-looking links. They seem perfectly normal until you click, and suddenly you're staring at a page just waiting to steal your bank details, passwords, or whatever else you meant to keep to yourself. Most folks can't tell the difference between a safe link and a scam, and that's really why online fraud keeps going up. That's exactly why we're building this project. It's a machine learning tool that checks sketchy links before you get burned. It doesn't just skim the URL—this thing looks at the length, weird symbols, sketchy domain names, HTTPS or not, scammy-sounding keywords, the whole lot. We train the model to find all those trouble signs, then it puts each URL into buckets like safe, phishing attempt, or defacement. The tool itself is just a web app. Drop in a link, hit check, and in seconds you get an answer: good to go, or risky. You even see a risk score, so you know if it's best to run the other way. If it's especially bad, you get an alert on the spot. On top of that, there's more—like digging into a site's background, watching for strange patterns, and tracking URLs in real time, so nothing slips by. It's all run through an easy dashboard, so you see the latest scans, what threats are out there, and anything new the system flags. With machine learning and real-time tracking, it's not just waiting for trouble—it actually goes out and finds it before it hits

Keywords —Machine Learning, Malicious URL Detection, Phishing Detection, Web Security, Cyber Threat Intelligence, URL Feature Analysis, Risk Score Assessment, AI-Based Cybersecurity, Suspicious Link Detection, Real-Time URL Monitoring, Web Threat Analysis, Intelligent Security Systems.

I. INTRODUCTION

The internet pretty much runs our lives these days—how we talk, shop, or dig up information. More people log on every year, and most just start banking on their phones, buying stuff, posting everything online, and tossing files into the cloud without a second thought. It's quick and easy, but honestly, that speed comes with a price: cybercriminals have a bigger target than ever. They're always out there, pulling new tricks. Phishing emails, bogus sites, fake login pages—they're just part of the scenery now. And it really does only take one click on the wrong

link to land yourself in trouble. Fake URLs blend right in these days. Most of us can't spot them, which makes it way too easy to hand over your password or bank details without even realizing. No wonder those scam stories keep popping up.

Old-school security used to fight back with blacklists and simple rules, but that only works against threats we already know about. Hackers don't hang around—they toss out new phishing sites, mess with domains, use link shorteners, whatever it takes to slip past basic defenses. Those simple roadblocks just can't

keep up. What we really need are smarter, faster ways to catch shady links—something that adapts as quickly as hackers do.

That's where machine learning and AI step in. They don't just scan for the usual red flags. Instead, they sift through mountains of data, spotting anything out of place. They look at link length, odd spellings, weird symbols, suspicious keywords, even fake security signals. Feed them thousands of real and fake links, and they actually get good at seeing the tricks—catching bad sites and saving you before you click.

That's what this project does. It's a sharp URL checker fueled by machine learning. Paste in any link, and it breaks it down, runs it through a trained model, and spits out an answer in seconds: harmless, phishing, or defaced. And it all lives in a web app—just drop in your link, get a result. Sometimes that quick pause is all you need to avoid disaster.

But there's more to it. Instead of just giving you a basic yes or no, the system gives you a risk score. You really see how shady (or safe) a link looks before you do anything. There are bonus checks, too—like scanning the domain's reputation, spotting weird behavior, and tracking changes as they happen. All these parts come together to spot threats early and catch stuff that basic tools would miss.

If you're a security pro or just someone who wants to stay out of trouble, you get an easy dashboard with everything laid out: how many links got checked, which threats showed up, what's still risky, and what's new. It's simple and fast. You get answers right now, so it's a lot easier to tell if a link is really safe.

Bottom line? Mixing machine learning with real-time checks actually gives us a shot at keeping up. Old tools chase after threats, but this new approach keeps learning and stays a step ahead. As the web gets trickier, tools like this go from nice-to-have to totally essential. Smarter link protection just makes the internet safer for everyone.

II. ALGORITHMS

The Machine Learning-Based Suspicious

URL Detection System is used to identify harmful contents within the internet using algorithms. In a world where numerous people are trying to trick people to visit their sites with fake websites and phishing attacks, it's necessary to have a system that can scan a URL to determine whether or not it is safe to proceed to. By using machine learning and specific algorithms, this system is able to identify certain patterns that are associated with each of the 3 types of URLs and classifies them into safe, malicious, or fake. By using algorithms, this system can recognize a wider range of patterns and is less susceptible to misclassifications compared to rule-based systems.

The system has a number of important algorithms within it that are responsible for the performance of the system. These algorithms are primarily used to analyze and identify the characteristics of a URL, to classify a URL based on machine learning principles, to calculate the risk associated with a URL, to detect threats in real-time, to monitor the threats found, and to avoid false detections. Here's an explanation of the main algorithms present within the system:

A. URL Feature Extraction Algorithm

The feature extraction algorithm is a critical algorithm in the process of detecting malicious URLs. It is responsible for analyzing the actual text of the URL and identifying its characteristics or 'features' that will later be used for classification. The system first preprocesses a URL when an analysis request is submitted to it by cleaning up the string and ensuring that the format is correct. After this it determines and records features such as the total length of the URL, the number of dots and special characters present within it, the total number of sub-domains, if there are IP addresses present within it, if the URL uses HTTPS, and if certain words like 'login', 'verify', etc, are present within it.

The system then uses these features to identify specific patterns that are common in most malicious websites. For example, the features in a phishing URL often include numerous strings of characters and numbers as opposed to a real domain name, multiple dots and hyphens within the domain, and the presence of specific keywords. The system will convert all identified features into an array of integers which will then be used by the machine learning algorithm for classification. This

algorithm enables the system to analyze and learn the structure of URLs to detect unusual and malicious patterns.

B.Machine Learning-Based URL Classification Algorithm

The most essential component of the system is the machine learning-based URL classification algorithm. Its sole function is to classify a given URL as either good, malicious or fake. The system employs a technique called machine learning or training where a pre-trained model uses existing labeled URLs (known good and bad URLs) as data to learn patterns that help distinguish the two types of URLs. During this training, algorithms such as Decision Tree, Random Forest, and XGBoost examine the features of each URL and learn to associate them with their correct classification.

When a URL is submitted for analysis, the system will examine its features and then use the trained model to predict the type of URL it is likely to be. Each prediction is assigned a score for each category (good, malicious and fake) and the URL is ultimately classified as the type with the highest score. By using the machine learning classification algorithm, the system is able to achieve high accuracy in determining the nature of the URL as compared to systems that use rule-based detection methods.

C.Risk Score Calculation Algorithm

Besides identifying whether a URL is bad or not, the system can also determine how risky it is for a user to access a given URL by assigning a risk score. This risk score is calculated by taking into account a multitude of factors which includes how confident the machine learning model is with its prediction, the presence of numerous bad words within the URL, if the URL has a strange domain name, and other such features that indicate how malicious a URL could be. All these individual components are given a weight to help determine the overall risk score of the URL.

The risk score ranges from 0 to 100 where a score of 0 indicates a perfectly safe URL while a score of 100 signifies a highly malicious URL. Based on this score, the system can categorize the URL into several levels such as 'safe', 'suspicious', and 'high risk'. When the risk score is too high, the system would notify the user about how dangerous

the URL is, compared to just knowing it is good or bad.

D.Real-Time URL Analysis Algorithm

The real-time analysis algorithm allows the system to check a URL immediately when a user enters it into the web application. When the system is given a URL to analyze it would immediately scan and identify its features, feed them into the machine learning algorithm, and compute a risk score all at once so that the user is made aware of how safe it is for them to access that URL before they attempt to visit it. This algorithm is particularly useful in preventing users from accessing malicious sites thereby increasing the overall effectiveness of the system's cyber security functions.

D.Threat Logging Algorithm

The threat monitoring algorithm helps keep a track of all the URLs scanned and the classification assigned to it along with its risk score. It would also maintain records of all the types of threats detected (good, malicious and fake) with their assigned risk score. This algorithm enables the system to have a record of all the URLs previously analyzed and it is useful for identifying particular types of malicious websites and online threats by monitoring them.

Additionally, the threat monitoring algorithm presents information regarding the number of bad URLs found on the dashboard, the types of phishing attacks that are most prevalent, and how often malicious URLs are being detected on the web, all of which assists in administering the system and spotting security threats effectively.

E.False Detection Reduction Algorithm

To further improve the reliability of the system, a false detection reduction algorithm is implemented. This algorithm re-examines the score given by the machine learning classification model, and combines the results of other detectors (such as keyword detection or domain pattern detection) to confirm whether a particular URL is genuinely malicious or not. If the model provides a high probability for it being malicious but the other detectors show conflicting information, it's important to recheck as the model could be giving a false positive. By performing further checks, this algorithm makes sure that not all good URLs are

flagged as bad but all bad URLs are identified effectively.

III. PROPOSED SYSTEM

The proposed system is an intelligent URL security analysis platform for the detection of malicious and suspicious URLs by using machine learning algorithms. Since there is an increase in phishing attacks, malicious websites, and web defacement incidents, it has become very important for users to get any tool that can analyze and check the security of the given URLs automatically before they navigate them. The proposed system uses the various features of the URLs from structural, and lexical features and classifies the URLs into Benign, Phishing, Defacement, etc based on these features by using ML techniques. System identifies the malicious and dangerous links by using various URL structures, special characters, domain structures, suspicious keywords, HTTPS usage etc.

The proposed system comprises feature extraction algorithms, ML classification models, and a risk scoring mechanism. After the user submits the URL, the features of the URL are automatically extracted and sent to the ML classification model that is trained by ML models like Random Forest or XGBoost. ML model predicts that the given URL is either a malicious or a benign link by analyzing all the features. Not only the classification, it gives a risk score that specifies how much potential danger the URL poses. If the value of the risk score is more than a predetermined value, the system gives an immediate warning that the URL is dangerous.

The system has a feature of real-time monitoring and analysis which helps to prevent the users from accessing the suspicious links before interaction. By analyzing the malicious links it will save users from data theft, phishing and malware infections and gives an intelligent and effective URL security analysis system for the users.

IV. WEB AND MOBILE PLATFORM BASED

The proposed URL detection system is implemented on a web platform, so the users can easily get the malicious URLs checked via an

interactive and real-time interface. After giving the URL as an input, the system immediately processes the URL using the feature extraction algorithm, machine learning classification algorithm, and a risk scoring algorithm, and provides the user with the result of whether it's a malicious or benign link and its risk score and threat category respectively.

The platform provides a dashboard to the user which contains the visualization and statistics of URL analysis. The dashboard shows the total number of URLs scanned, number of phishing attacks detected and URL analysis results. Using this dashboard, the user or administrators can have a track of the overall system activity and the suspicious URLs.

The backend system of this proposed URL security analysis system is built by using the programming language Python and the framework Flask. The ML model is trained by using various Python libraries like Scikit-learn, Pandas, and NumPy, etc. The training model is saved in a serialized format, so that it can be accessed quickly for predictions in real-time on the web platform and increases the efficiency, scalability and accuracy of the system.

V. AI-POWERED URL ANALYSIS AND THREAT DETECTION

For cybersecurity professionals, minimizing false positives when analyzing malicious URLs is essential. The URL detection system has been designed to address this with sophisticated pattern analysis and ML algorithms. After a URL has been submitted by the user, features are extracted by the application in order to parse the structural and lexical components of the URL. These components include: URL length, the amount of special characters in the URL, the amount of sub-domains in the URL, the domain structure of the URL, if IP addresses are used in the URL, and common keywords present in the URL (such as "login," "verify," and "update") and if the URL uses an HTTPS protocol.

These features are then fed into a supervised ML algorithm such as Random Forest or XGBoost that has been trained on a database of URLs (including benign, phishing and defaced URLs). Through this model, a ML algorithm

learns these characteristics and is able to predict whether a submitted URL is a safe or malicious one. Threat level and risk score is determined by the probability that the ML model's prediction will be accurate, thus determining if the URL is a medium, high or low security threat. This AI model offers better efficiency due to its ability to analyze future attacks of an unknown nature based on established patterns discovered through the large quantities of data.

VI. INTELLIGENT URL ALERT SYSTEM

The intelligent alerting function of the system will inform the user when they attempt to access a potentially dangerous URL. When a threat level has been identified the system immediately warns the user if a URL has been identified as malicious or defaced. The alert will provide the predicted class and risk score of the URL along with the reason why the specific URL is suspect.

A database for threat monitoring can also be developed where recurring patterns in the structure of a malicious URL can be recorded in order to monitor activities by administrators and types of cyber-attacks. In this manner users are alerted prior to visiting a website that may potentially harm them.

VII. SECURE DATA MANAGEMENT AND PRIVACY

Data security and privacy have always been crucial to an application dealing with potentially sensitive data. The entire system has been designed around storing user input data and results data securely. All logged information including analysis results and operation logs is encrypted with AES prior to storage within the database.

Data and operation logs are stored on large data stores such as those provided by Firebase or AWS in the case of high volume operations and limited system access rights is enforced by the administration in order to maintain user privacy and prevent any confidential information from leaking.

VIII. REAL-TIME URL ANALYSIS AND MONITORING

The real-time URL analysis gives the user

an instantaneous reply with regards to how safe the submitted URL is. The features of the URL are extracted by the application, transmitted to the ML algorithm located on the backend server, where the prediction of if the URL is safe or not takes place instantaneously. This classification is then passed back to the user within a few seconds time so that they can immediately determine whether or not to continue to the page.

The system also allows for the monitoring function to take place where all URLs, the classification of those URLs and their associated threat level are logged for analytical purposes. All logged information is displayed on the dashboard where you can then see statistics including how many URLs have been scanned, the distribution of threat levels amongst URLs submitted, and recent activity, in order to keep administrators well informed about the nature of current cyber-threats.

IX. USER AWARENESS AND CYBERSECURITY EDUCATION

The purpose of the system also extends beyond merely detecting malicious URLs, it also aims to improve user awareness of cybersecurity. Users are given safety warnings and education as well as being given warnings in the case that they submit a suspicious link. Phishing, domain names and basic web browsing habits are explained to users in order to improve their understanding of cyber-attacks. With the right education and awareness of potential threats it is more likely that the user will practice safe web browsing and fall victim to fewer phishing attacks.

X. TECHNOLOGY

Machine learning and web technologies along with secure computing paradigms are utilized in this URL detection system. Algorithms like Random Forest or XGBoost are used to classify the URL. The application is programmed in Python on the back-end, using Python libraries such as Scikit-learn, Pandas, and NumPy for the ML components. The web part is built using web technologies (HTML, CSS and JavaScript) to construct the interface. A stored serialized object will hold the trained ML algorithm so that it can be used by the web application, whereas system logs can be transmitted to a cloud data storage service (e.g. Firebase or AWS) where data can be

managed and monitored in a scalable fashion.

Through utilizing both machine learning techniques and real time data analysis with careful management and monitoring of information, an intelligent system is built that can identify malicious URLs and defend users in the cyber world.

XI. EXPECTED OUTCOMES

The designed Machine Learning-Based URL Security Detection System has various advantages when it comes to online security and awareness. The system is able to give users instant confirmation of the safety of any given URL before they access it and thus help them avoid falling into phishing websites, infected with malware and fraud. In addition to being based on ML algorithms, the system also takes advantage of advanced feature extraction to accurately identify malicious URL patterns into either benign, phishing or defaced categories. The introduction of a risk scoring system for every URL alerts users of the threat level and allows them to decide wisely about whether to proceed or not while navigating online.

With a real time analysis capability the system is also able to identify dangerous links instantly while the monitoring dashboard can give administrators an overview of the URL scan statistics and upcoming cyber threats. This ML system also makes users aware of any suspicious URL patterns that should be avoided and ultimately serves to achieve digital security through accurate threat detection using the above mentioned approach.

XII. RESULTS AND DISCUSSION

The construction and operation of a Machine Learning-Based Suspicious URL Detection System is considered a valuable advancement in the pursuit of cyber security and the prevention of threats related to malicious websites. By incorporating both machine learning and feature extraction methods along with web technologies, a number of URLs are analyzed for threats. In order to estimate the accuracy and performance of the implemented detection model, experimental testing was carried out using a dataset comprising malicious, phishing and defaced URLs.

Through testing, multiple features of each URL including length, the inclusion of special characters, the arrangement of domains, suspicious key words and whether or not HTPSS is being used, were processed through the developed ML model which was trained accordingly to determine malicious or normal patterns associated with certain URLs. Results obtained from the test prove that the model can effectively predict suspicious features prevalent in phishing URLs or otherwise. The system generated a risk score corresponding to the degree of threat for each URL so that users are appropriately alerted to potentially risky links or websites and their probability of being dangerous is shown as 'risk score'. The system also has a web based interface for real time detection so that URL checking occurs at the very moment a user chooses to navigate to the said link.

The dashboard section for the ML based system aids in the visualization of these statistical data collected on the URLs scanned, as well as giving a bird's eye view of cyber threats which may be prevalent. From this the user or administrator may ascertain the threat rate of web links that may cause a problem for the user. The outcomes from the test clearly show the effectiveness of the developed system in detecting threats so that users may be immediately warned, thus preventing them from being victims to such attacks.

XIII. CONCLUSION

The development of the Machine Learning-Based Suspicious URL Detection System offers an intelligent and effective approach for preventing users from falling to harmful web links, and ensuring robust cybersecurity. In comparison to traditional detection methods that largely use blacklists that are not able to recognize new phishing websites, this ML system is able to achieve its task by observing patterns in URLs.

It can effectively detect phishing sites, dangerous URLs and malicious web defacements through the use of feature extraction algorithms and ML classification models to create accurate threat identification through real-time analysis. The use of the risk scoring system is user friendly and is informative about the level of threat from any given URL. It is easy to implement with its

web interface and the dashboard allows for users and administrators to monitor statistical data and threats. This project shows that ML can be used to great effect in making cybersecurity systems stronger, not just in terms of detection, but also user awareness. If developed with more advanced ML models and external threat intelligence sources, this system can evolve into a robust security platform for protection against the ever evolving cyber threats.

REFERENCES

- [1] M. Verma and N. Kumar, "Detection of Malicious URLs Using Machine Learning Techniques," *International Journal of Computer Applications*, vol. 178, no. 40, pp. 15–20, 2019.
- [2] J. Ma, L. K. Saul, S. Savage, and G. M. Voelker, "Beyond Blacklists: Learning to Detect Malicious Web Sites from Suspicious URLs," *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1245–1254, 2009.
- [3] A. K. Singh and B. B. Gupta, "Phishing Detection Using Machine Learning Techniques," *International Journal of Information Security*, vol. 19, no. 3, pp. 287–299, 2020.
- [4] S. Sahoo, B. Liu, and S. C. Hoi, "Malicious URL Detection Using Machine Learning: A Survey," *ACM Computing Surveys*, vol. 52, no. 1, pp. 1–36, 2019.
- [5] Y. LeCun, Y. Bengio, and G. Hinton, "Deep Learning for Cybersecurity Applications," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [6] R. Purushothaman, A. K. Thejasree, and K. Udaya Bhaskar, "Secure Web and URL Analysis Using Machine Learning," *International Journal of Research in Engineering, Science and Management*, vol. 5, no. 4, pp. 230–235, 2022.
- [7] A. Chaudhari and P. R. Bhaladhare, "A Survey on Machine Learning Techniques for Detecting Phishing Websites," *International Journal of Advanced Research in Computer Science and Software Engineering*, vol. 8, no. 6, pp. 45–52, 2021.
- [8] S. B. Rathod and R. P. Patil, "Detection of Phishing Websites Using Random Forest Classifier," *International Journal of Computer Science and Information Technologies*, vol. 11, no. 2, pp. 89–95, 2020.
- [9] T. Fette, N. Sadeh, and A. Tomasic, "Learning to Detect Phishing Emails and Malicious URLs," *Proceedings of the International World Wide Web Conference*, pp. 649–656, 2007.
- [10] J. Zhang and C. Gupta, "Real-Time Phishing URL Detection Using Machine Learning Algorithms," *IEEE International Conference on Information Security and Privacy*, pp. 210–215, 2021.