



# INTERNATIONAL JOURNAL OF ADVANCE RESEARCH, IDEAS AND INNOVATIONS IN TECHNOLOGY

ISSN: 2454-132X

Impact Factor: 6.078

(Volume 7, Issue 2 - V7I2-1201)

Available online at: <https://www.ijariit.com>

## Detection and prevention of phishing websites using web crawlers

Poojhaa S

[poojhaakumar173@gmail.com](mailto:poojhaakumar173@gmail.com)

PSG Polytechnic College, Coimbatore,  
Tamil Nadu

Narmatha M.

[aishu7801@gmail.com](mailto:aishu7801@gmail.com)

PSG Polytechnic College, Coimbatore,  
Tamil Nadu

Sowmiya R.

[sowmiyaramasamy0002@gmail.com](mailto:sowmiyaramasamy0002@gmail.com)

PSG Polytechnic College, Coimbatore,  
Tamil Nadu

### ABSTRACT

*Phishing attack is the simplest way to obtain sensitivity information from innocent users. The aim of the phishers is to acquire critical data like user name, password and bank account details. The proposed system mainly focuses on detecting and preventing the phishing websites. A Web crawler is used to detect the websites. Detected phishing sites are added to the black list. The black list contains only fake websites. Web Crawling focuses on obtaining the links of the web page. Usually, a phishing website can be easily identified by its URL by its HTML code. Here the check website page makes users to be aware of phishing website and avoid users to become victim of such attacks. This software is very helpful for users to identify and prevent the phishing websites. For this purpose, we are developing a phishing website detection system.*

**Keywords:** Phishing, Web-crawlers, Phishing Attack

### 1. INTRODUCTION

Phishing is a new word produced from 'fishing', it refers to the act that the attackers allure users to visit a faked Web sites. It is a new type of network attack; the attacker copies the content from the website of a well-known company or a bank and creates a phishing website. The attacker keeps a visual similarity of the phishing website similar to the corresponding legitimate website to attract more users. The phishers must duplicate the content of the target site and they must use tools to (automatically) download the Web pages from the target site. The proposed system helps the user to detect and identify the phishing websites. It makes use of web crawler to crawl the hyperlinks in the website. The website does not contain any hyperlinks, and then the website will move to the black list. Those sites are fake and the user does not wish to access. The web crawler can verify the user's inputted URL; if the result is phished, then it warns the user to block the website by adding it to a black list. The user can scan website, if it is a phishing site, then it will give alert message for preventing the user.

### 2. LITERATURE SURVEY

#### A. Antiphishing to Protect Users from Phishing

AntiPhish is used to prevent users from using fraudulent websites which, in turn, could lead to a phishing attack. Here, AntiPhish traces the sensitive information to be entered by the user and alerts the user whenever he/she attempts to share his/her information on an untrusted website. The most effective explanation for this is to encourage users to approach trusted websites only. This approach, however, is unrealistic. In any case, the user can get tricked. It is therefore obligatory for the associates to present such explanations in order to overcome the problem of phishing. Widely accepted alternatives for the identification of "clones" and maintenance of phishing website records in the hit list are based on creepy websites.

#### B. Learning to Detect Emails from Phishing

The required process of reliability of the system on a trait intended to reflect the besieged deception of users by means of electronic communication is an alternative to detecting these attacks. This method can be used to detect phishing websites or text messages sent through emails used to trap victims. Around 800 phishing mails and 7,000 non-phishing mails are traced to date and over 95 percent of them are correctly identified on the basis of 0.09 percent of real emails along with the categorization. We should simply conclude with the methods of detecting deception, along with the changing nature of attacks.

#### C. E-banking phishing detection method using fuzzy data mining

Identifying and classifying phishing websites, primarily used for e-banking services, is very complex and dynamic. Some critical data mining techniques can prove an effective way to keep e-commerce websites secure because of the presence of different ambiguities in the identification, as it deals with the consideration of different quality variables rather than precise values. An efficient approach to resolve the "fuzziness" in the evaluation of the e-banking phishing website is used in this paper to detect e-banking phishing websites using an intelligent, resilient and successful model. To consider different successful

factors of the e-banking phishing website, the implemented model is based on fuzzy logics along with data mining algorithms.

**D. Quick Flux Phishing Domains Collaborative Detection**

Two methods are listed here to find the connection of evidence from multiple DNS servers and multiple FF domain suspects. Examples of real life can be used to prove that our approaches to correlation expedite the FF domain detection, which is based on an analytical model that can calculate different DNS queries used to verify an FF domain. It also demonstrates that using a distributed model, which is more modular compared to a centralized one, is to publish N subscribe correlation model known as LARSID, to introduce correlation schemes at an enormous level. It is very difficult to classify the FF domains in an accurate and timely manner in deduction, as the proxy screen is used to protect the FF Mother ship. To evaluate the problem of FF detection, a computational approach is used by measuring the number of DNS queries needed to get back a certain number of unique IP addresses.

**3. EXISTING SYSTEM**

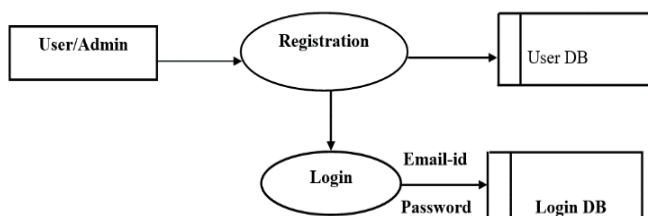
The Existing system determines manually whether the site is a phishing site or not, but it's difficult to find those phishing sites. The Web master of a legal Web site periodically scans the root DNS for suspicious sites (e.g. www.lcbc.com.cnvs.www.icbc.com.cn). Since the phisher must duplicate the content of the target site, he must use tools to (automatically) download the Web pages from the target site. It is therefore possible to detect this kind of download at the Web server and trace back to the phisher. For DNS scanning, it increases the overhead of the DNS systems and may cause problem for normal DNS queries, and furthermore, many phishing attacks simply do not require a DNS name. The clever phishers may easily write tools which can mimic the behaviour of human beings to defeat the detection.

**4. PROPOSED SYSTEM**

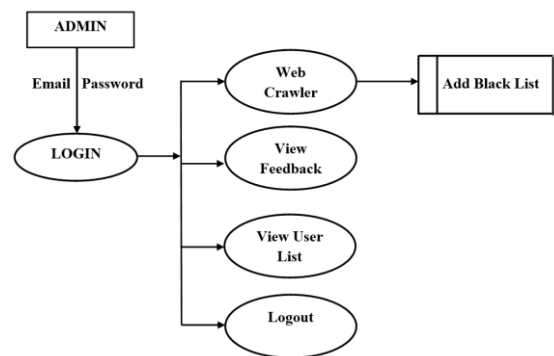
In proposed system to detect the phishing website using web crawler. The search starts by crawling the pages of your site. Then it continues to visit the links (web page addresses or URLs) that are found in your site. In our proposed system the admin can login and enter any website URL in web crawler, then it searches the URL, it identifies the hyperlinks of the page. The links are display on current page. The website does not contain any hyperlinks, and then the website will move to the black list. The black list contains phishing sites. Those sites are fake and the user does not wish to access. Now, the user can enter the URL of website within the check website page, then click scan button, the scanning process will start and to check whether the entered URL is in the blacklist or not. If the site is not in the black list further process will be preceded. Otherwise, if the URL is in the black list, it is considered as a phishing site, and it will give the alert message to the user.

**5. DATA FLOW DIAGRAMS**

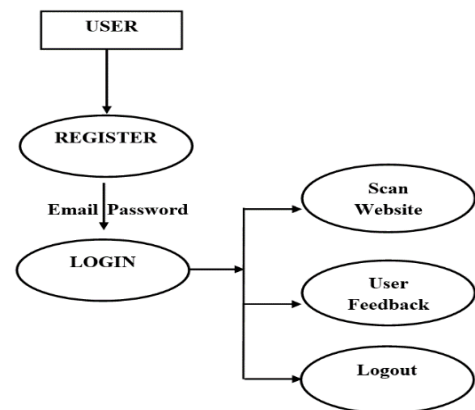
**LEVEL 0**



**LEVEL 1**



**LEVEL 2**



**6. MODULE DESCRIPTION**

**A. Login**

The Login page will be displayed to user when user enters after registering their details. Then user can enter an Email-id and Password. To start the process, the user has to login by clicking the login button.

**B. Registration**

In this registration module, the user has to give their details such as User name, Email-id, Password, Retype Password, Date of birth and Gender through the user can login. If the user is an existing user then they shall continue the process by giving their authentication details. If the user is not an existing user then they have to register by giving their details. MD5 algorithm is used to secure the user's profile by encrypting the password.

**C. Web Crawler**

Crawlers look at Web Pages and follow links on those pages. In this module we can give some website URL. If it is not a phishing site, it will crawl the links of that website; otherwise, that link will add to the blacklist.

**D. Blacklist**

The blacklist contains the phishing sites. Which works on list-based detection, known phishing sites which are in black list. And then these phishing sites are added to black list.

**E. Check Websites**

In this module user can enter the URL of website, to check if it is in the blacklist or not. If it is in the blacklist, those sites are fake and the user does not wish to access and then it will give the alert message to the user.

**F. Feedback Form**

User can give their feedback of the application. User can fill their details like Name, Email-id and then give the feedback.

**7. SYSTEM TESTING**

System testing is actually a series of different tests whose primary purpose is to fully exercise the computer-based system. Although each test has a different purpose, all work to verify that

all system elements have been properly integrated and perform allocated functions. During testing, it tried to make sure that the product does exactly what is supposed to do. Testing is the final verification and validation activity within the organization itself. In the testing stage, the following goals are achieved; it confirms the quality of the product, it finds and eliminate any residual errors from previous stages, it validates the software as a solution to the original problem, it demonstrates the presence of all specified functionality in the product, it estimates the operational reliability of the system.

## 8. SYSTEM IMPLEMENTATION

Implementation is the stage when the theoretical design is turned out into a working system. Thus, it can be considered to be the most critical stage in achieving a successful new system and in giving the user, confidence that the new system will work and be effective. The implementation stage involves careful planning, investigation of the existing system and its constraints on implementation, designing of methods to achieve changeover and evaluation of changeover methods. Implementation is the process of converting a new system design into operation. It is the phase that focuses on user training, site preparation and file conversion for installing a candidate system. The important factor that should be considered here is that the conversion should not disrupt the functioning of the organization.

## 9. RESULTS

Authorized users can access data by logging on this website as shown in Fig. 1. On entering the registered user ID and password, it goes to the webpage where the user should enter the URL of the webpage and enter the scan button. Once this is done the web crawler will scan the URL and notify the user whether the sight is phishing or not. So, this makes the user's details secured from the phishing website.

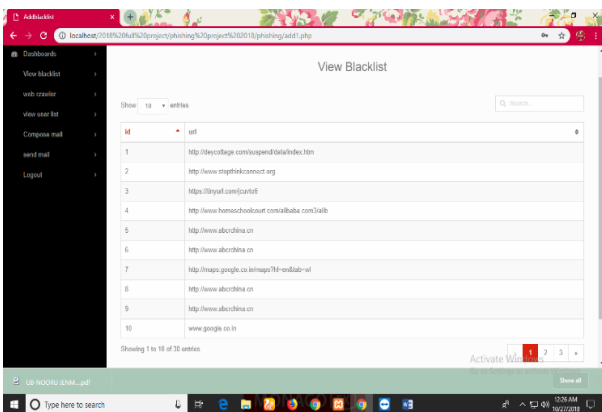


Fig. 1: View Black List

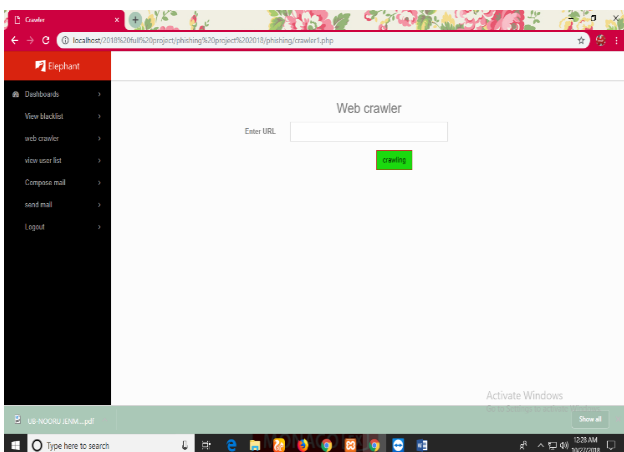


Fig. 2: Web crawler

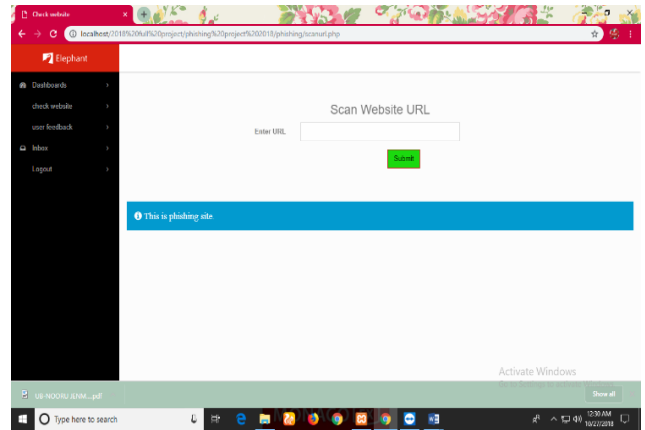


Fig. 3: Check Website

## 10. CONCLUSION AND FUTURE WORK

Phishing cannot be solved with a single solution. It is a critical situation in which Phishers always try to come up with brand new modes of manipulating the consumers. Online consumers should embrace regular risk scrutiny to detect the recent techniques which may head to a thriving Phishing attack. To find safer ways, user must be aware about the dangers of advanced malware which are taking place nowadays. Also, safekeeping teams need to execute advanced methodologies that can put the advanced threats to an end that are recently being bypassed by their predictable resentment. Further contribution is done in detecting the identity theft and the phishing mails. It does not involve in the rising trends towards e-mail outsourcing. Log analysis and communication taking place across managerial boundaries can prove to be a tricky one. In other words, we can also say that other electronic transactions will also become a part of the threats. Henceforth, it is suggested to sincerely work on these problems before attacks are being clutched wildly. A command should be acquired which can protect all crucial internet banking activities.

## 11. ACKNOWLEDGEMENT

We acknowledge with deep sense of reverence, our special gratitude towards our Head of the Department Dr.S.BRINDHA, Department of Computer Networking for her guidance, inspiration and suggestions in our quest for knowledge. We would like to express our gratitude towards our parents for their tremendous contribution in helping us reach this stage in our life. This would not have been possible without their unwavering and unselfish love, cooperation and encouragement given to us at all times. However, it would not have been possible without the kind support and help of many individuals. We would like to extend our sincere thanks to all of them.

## 12. REFERENCES

- [1] "Protecting Users Against Phishing Attacks with AntiPhish" Engin Kirda and Christopher Kruegel Technical University of Vienna
- [2] "Learning to Detect Phishing Emails" Ian Fette School of Computer Science Carnegie Mellon University Pittsburgh, PA, 15213, USA icf@cs.cmu.edu Norman Sadeh School of Computer Science Carnegie Mellon University Pittsburgh, PA, 15213, USA Anthony Tomasic School of Computer Science Carnegie Mellon University Pittsburgh, PA, 15213, USA
- [3] Phishtank - <https://www.phishtank.com/>
- [4] International Journal of Advanced Computer Technology (IJACT), "A Review of Various Techniques for Detection and Prevention for Phishing Attack".